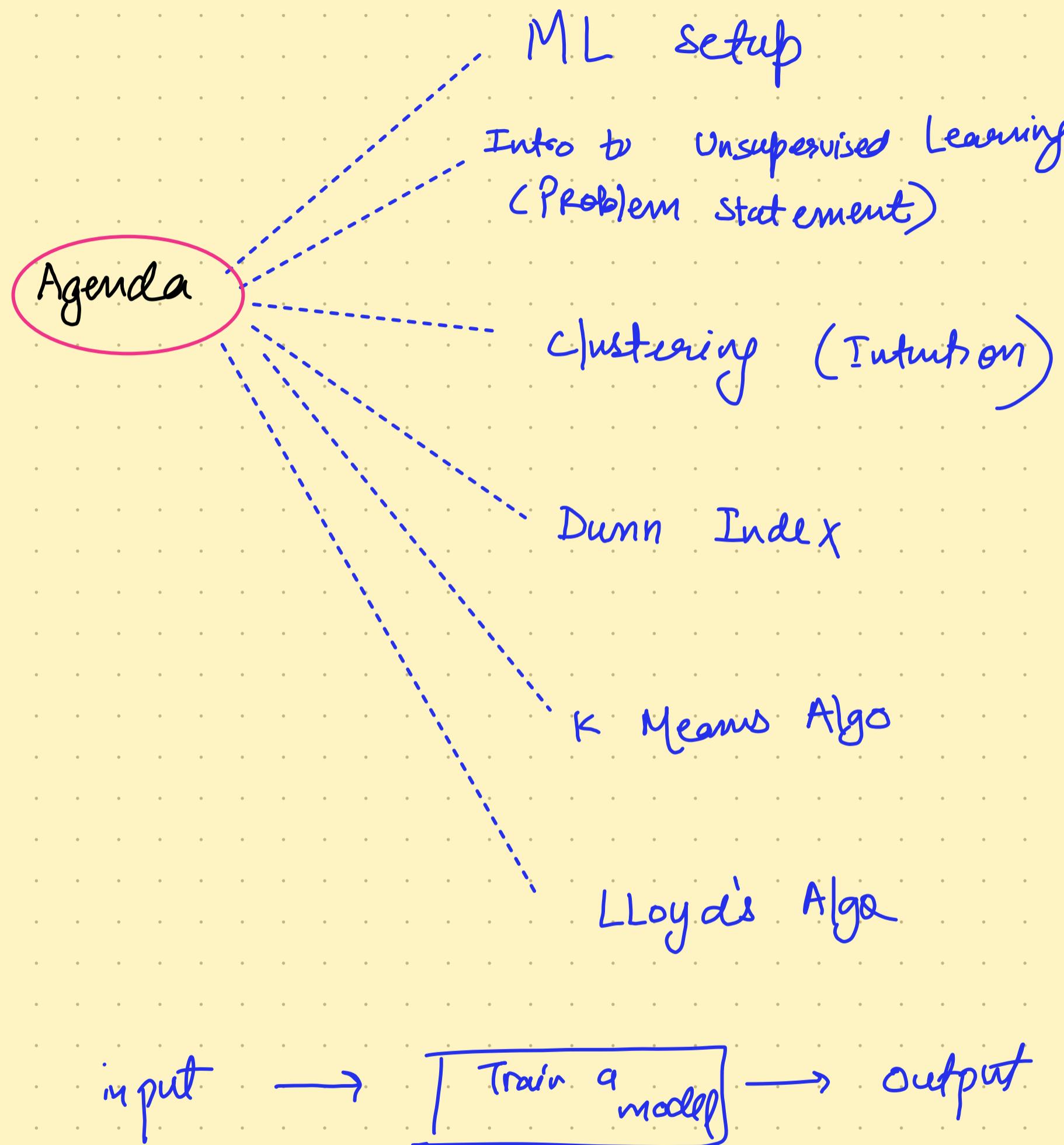


Class 1

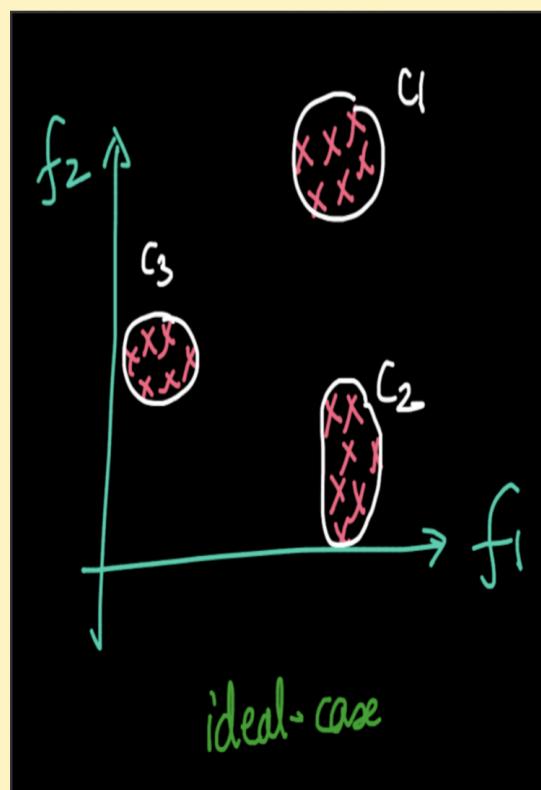
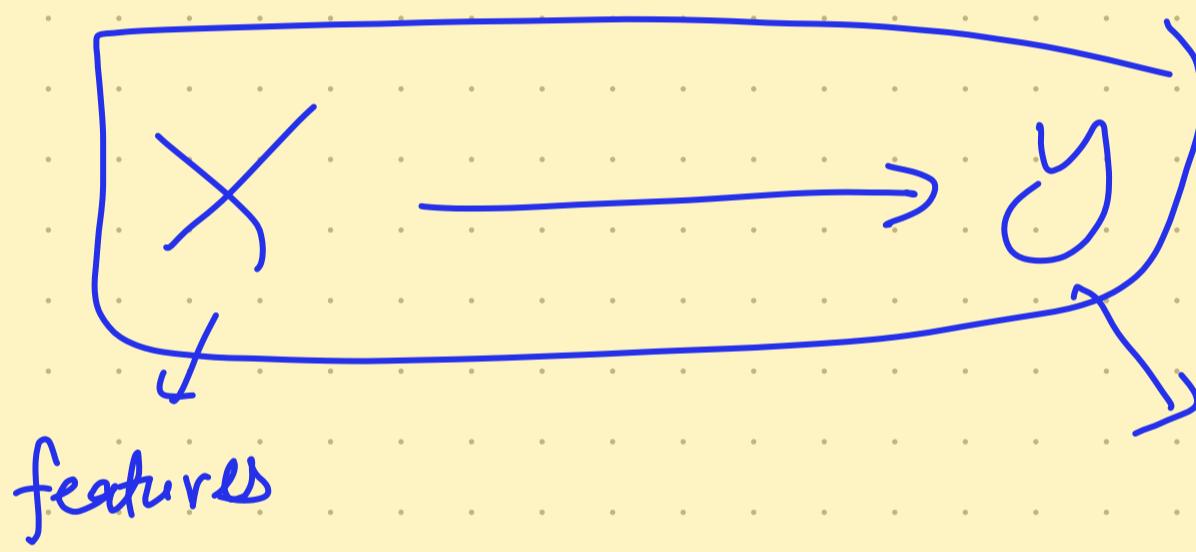
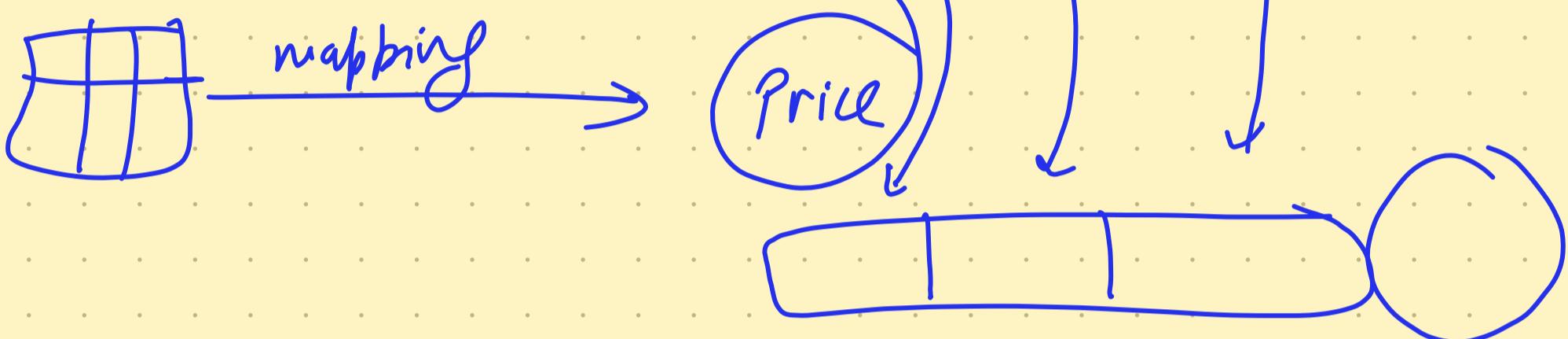
Introduction to Unsupervised Learning



Data

Training Data

Area m ²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000



INTUITION

- points in a cluster are close to each other
- points in different clusters are far from each other

Amazon Data

ID	n_clicks	n_visits	amount_spent	amount_discount	days_since_registration
0	1476	130	65	213.905831	31.600751
1	1535	543	46	639.223004	5.689175
2	1807	520	102	1157.402763	844.321606
3	1727	702	83	1195.903634	850.041757
4	1324	221	84	180.754616	64.283300

X labels

Supervised
M.L

House features

Price

House	Avg	No of bed room

unsupervised

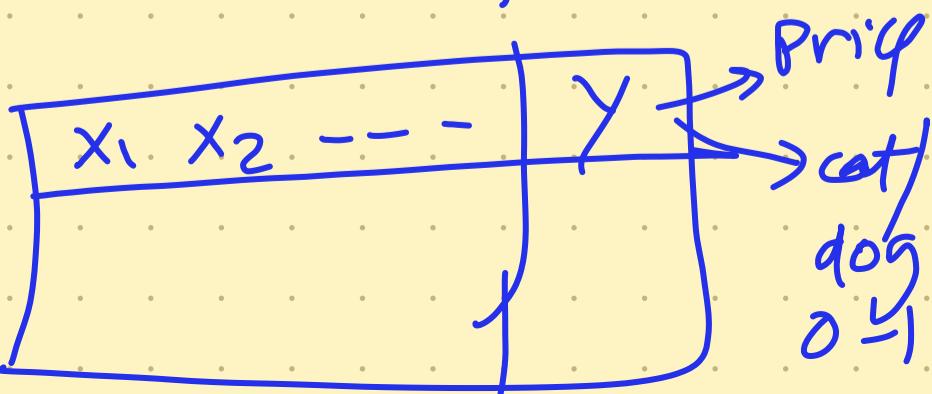
Supervised Machine Learning

Regression

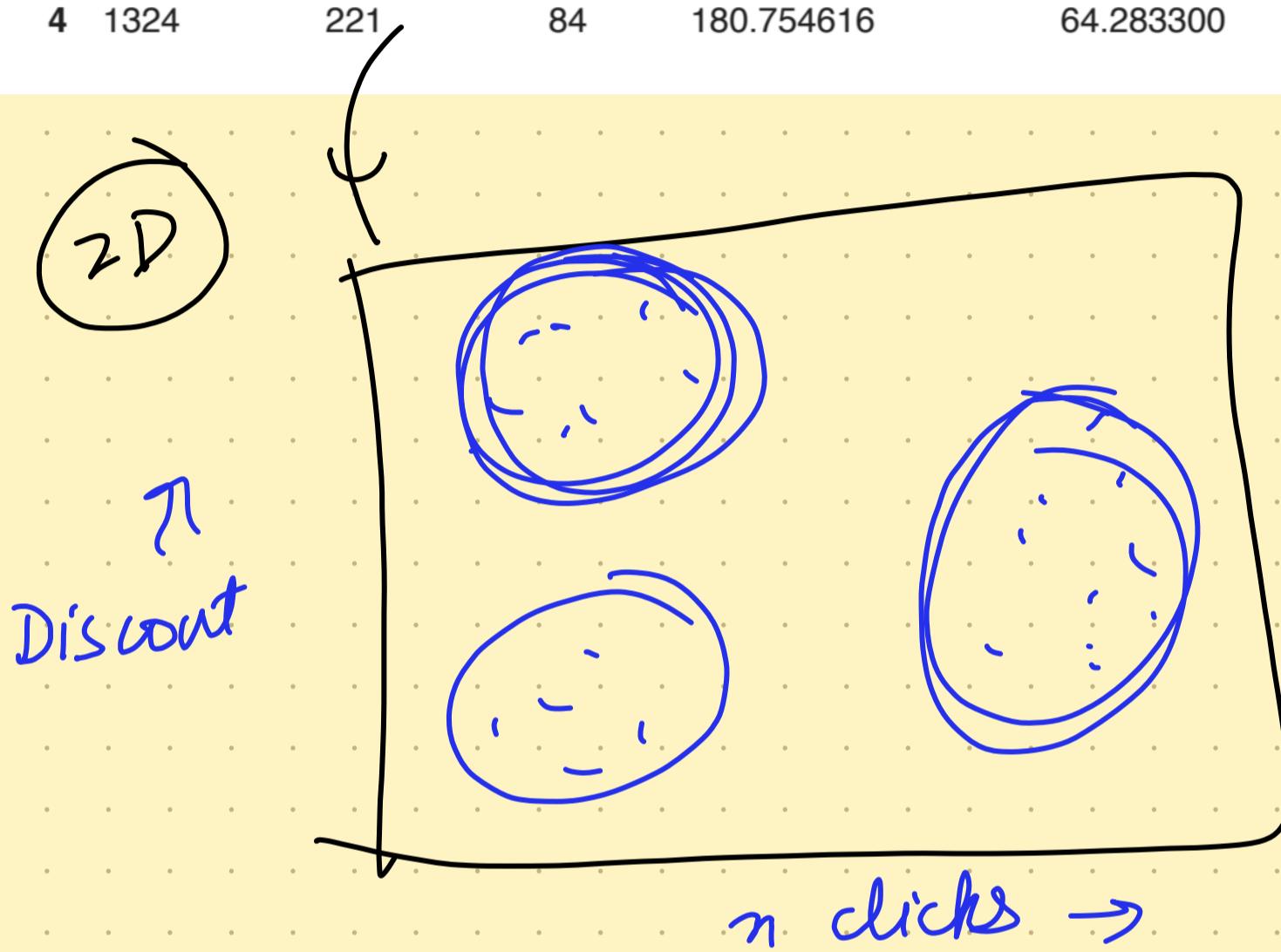
()

Classification

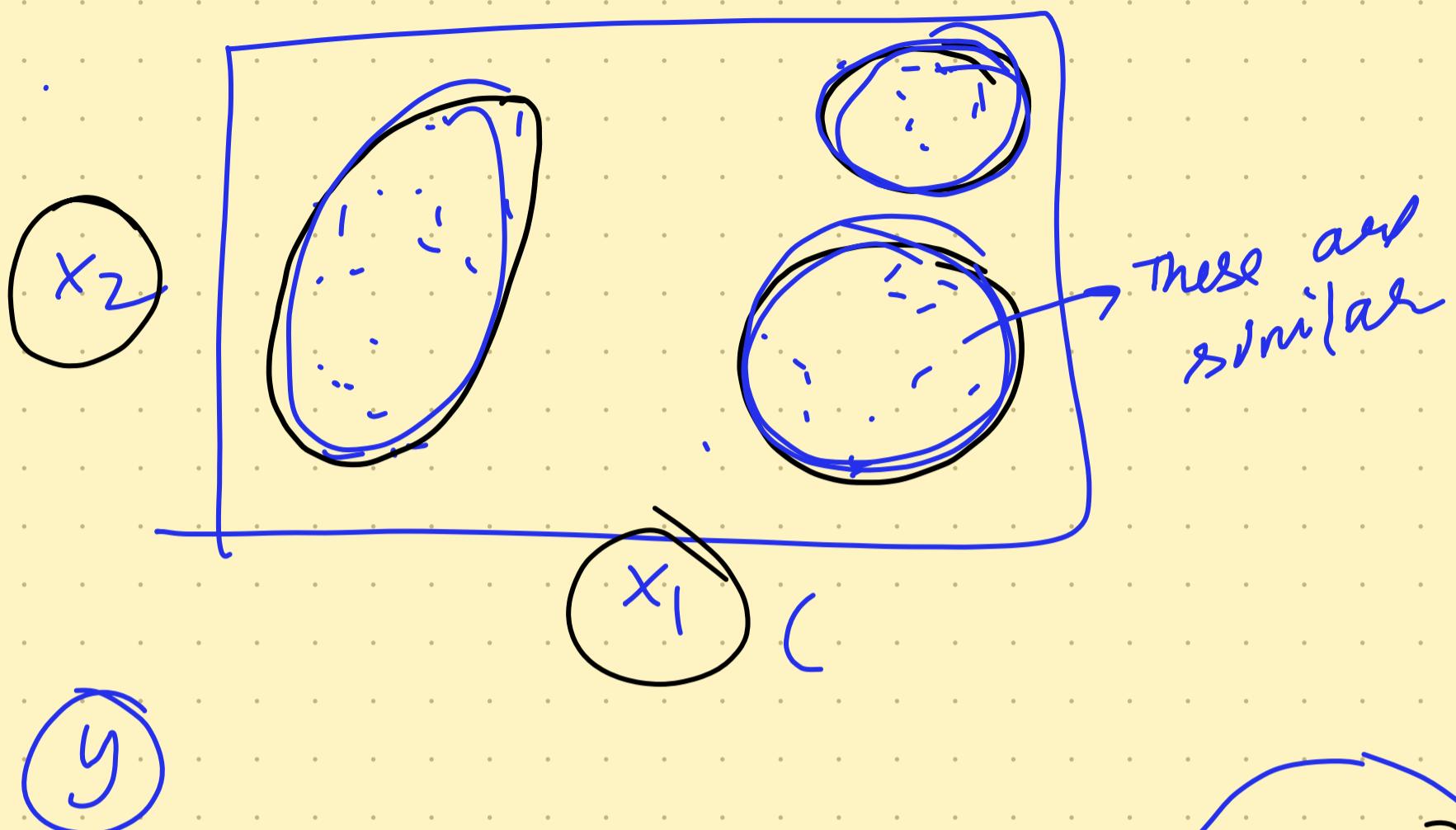
X and Y



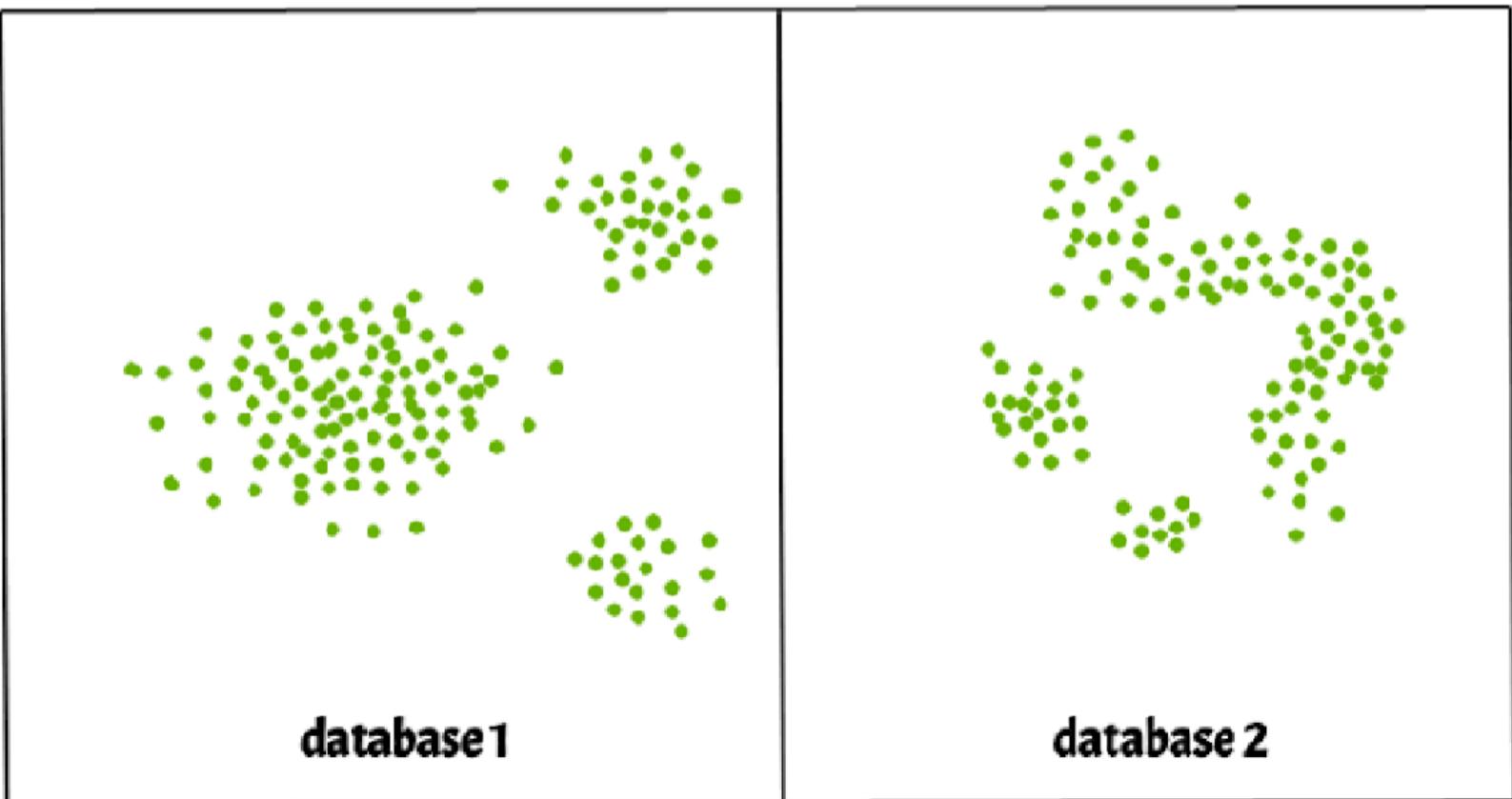
	ID	n_clicks	n_visits	amount_spent	amount_discount	days_since_registration
0	1476	130	65	213.905831		31.600751
1	1535	543	46	639.223004		5.689175
2	1807	520	102	1157.402763		844.321606
3	1727	702	83	1195.903634		850.041757
4	1324	221	84	180.754616		64.283300



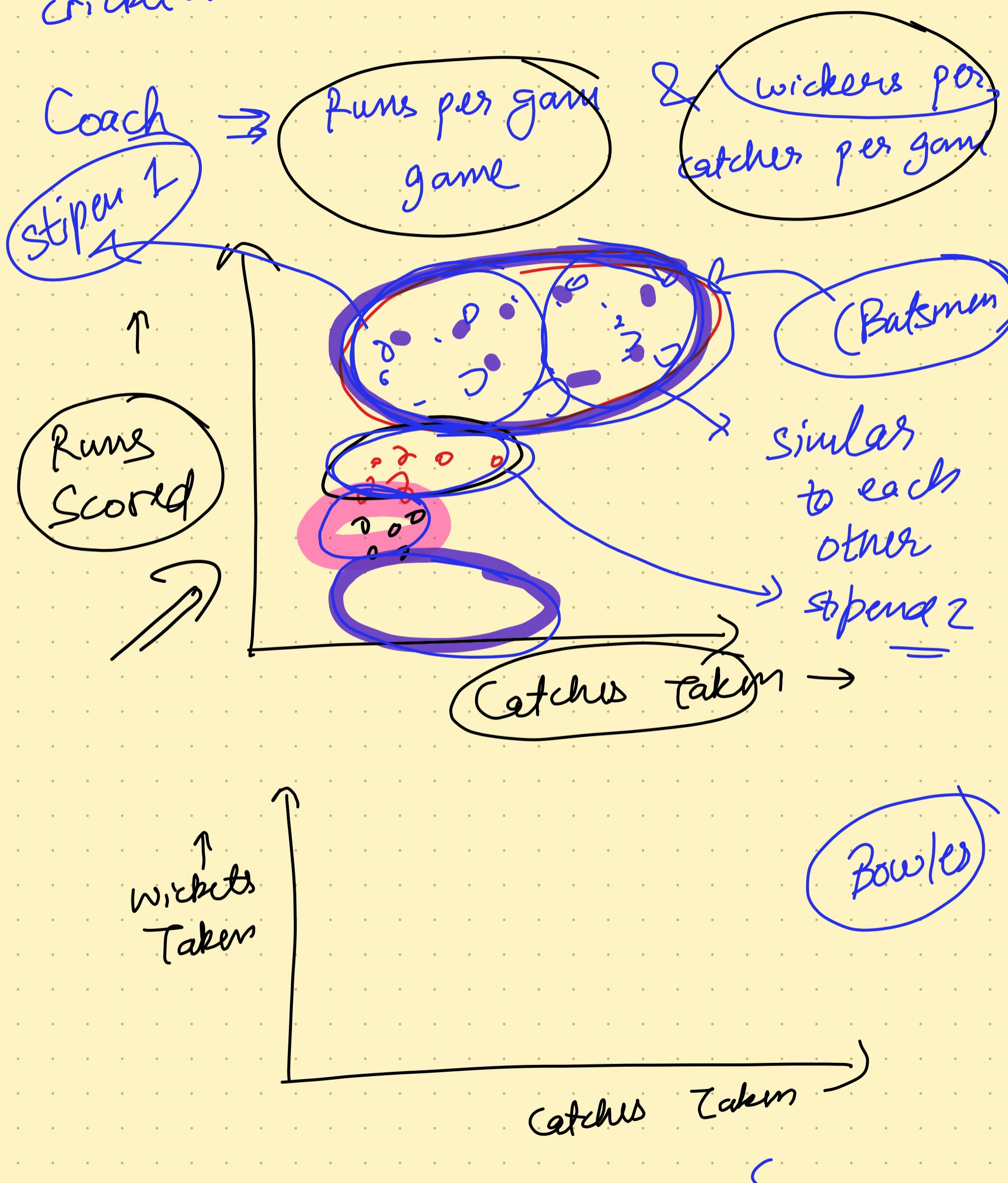
Unsupervised Learning →



Clustering [cricket Academy] → { 30,000 enrollment }

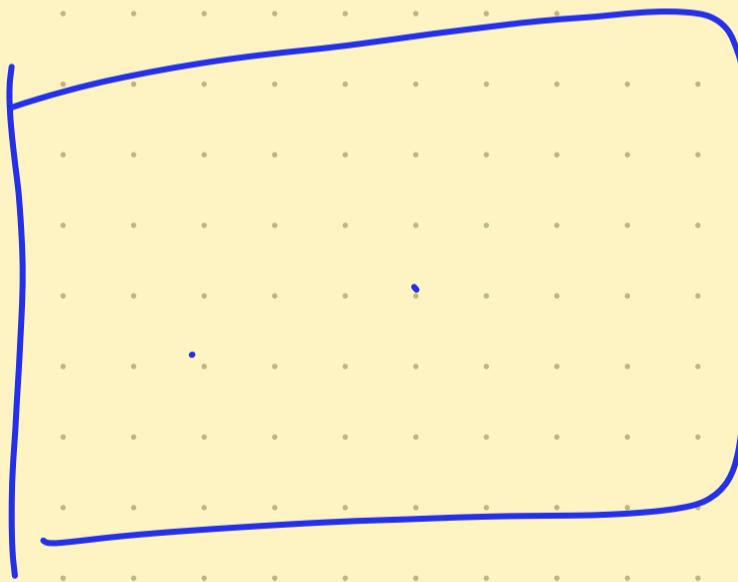


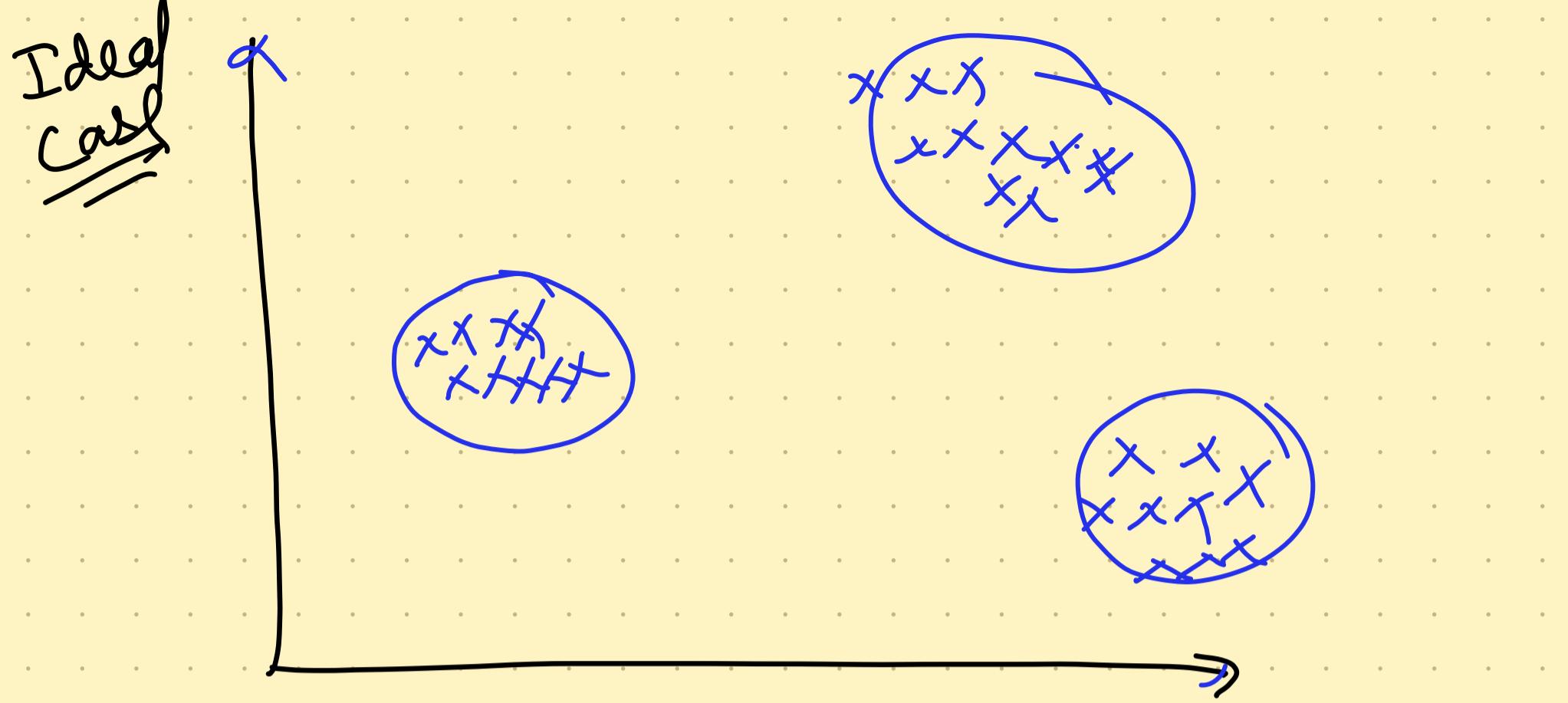
Stipend → Based on the performance of
cricketers



Player	Catches	Runs Scored
1	-	-
2	-	-
3	-	-
:	:	:

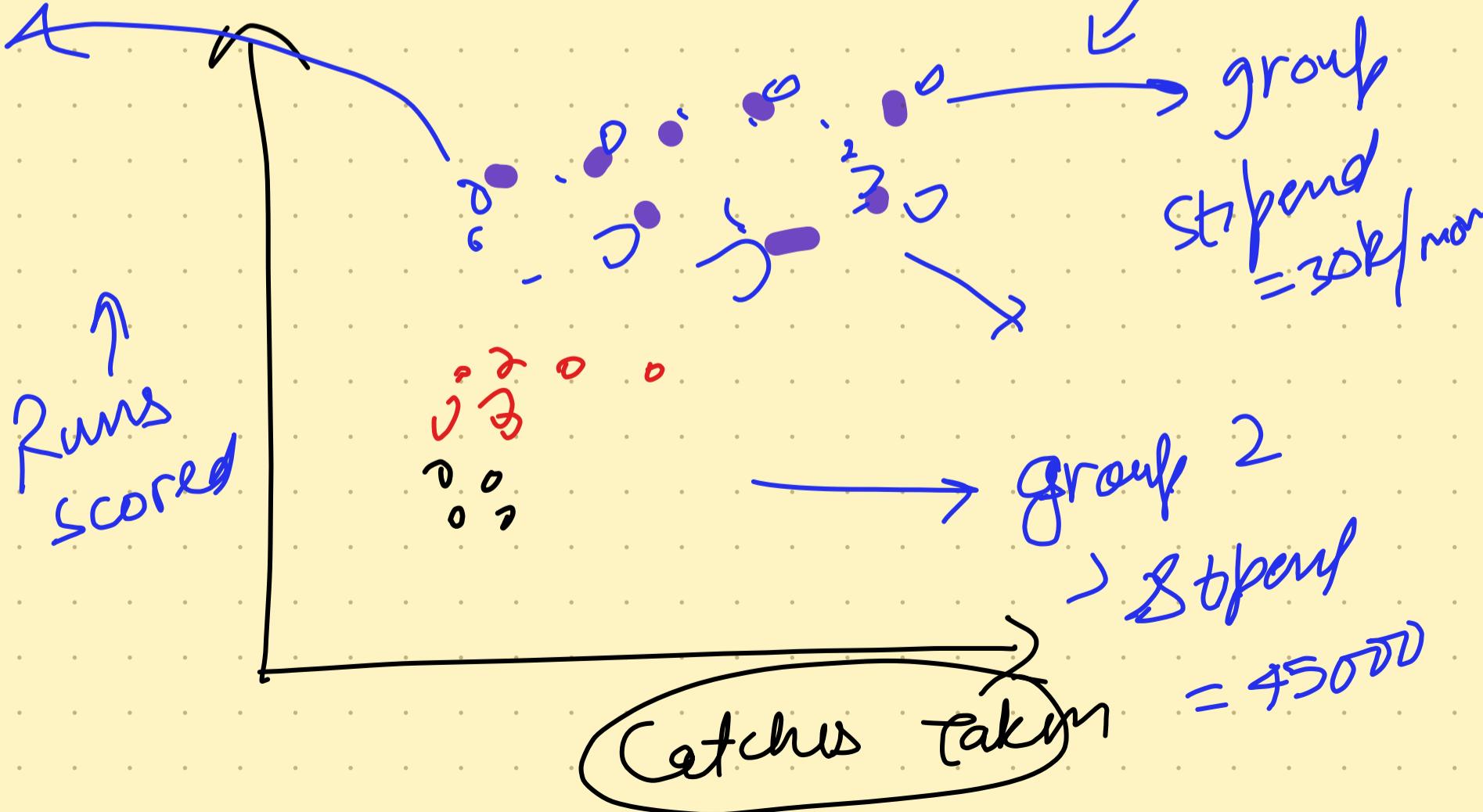
Isod

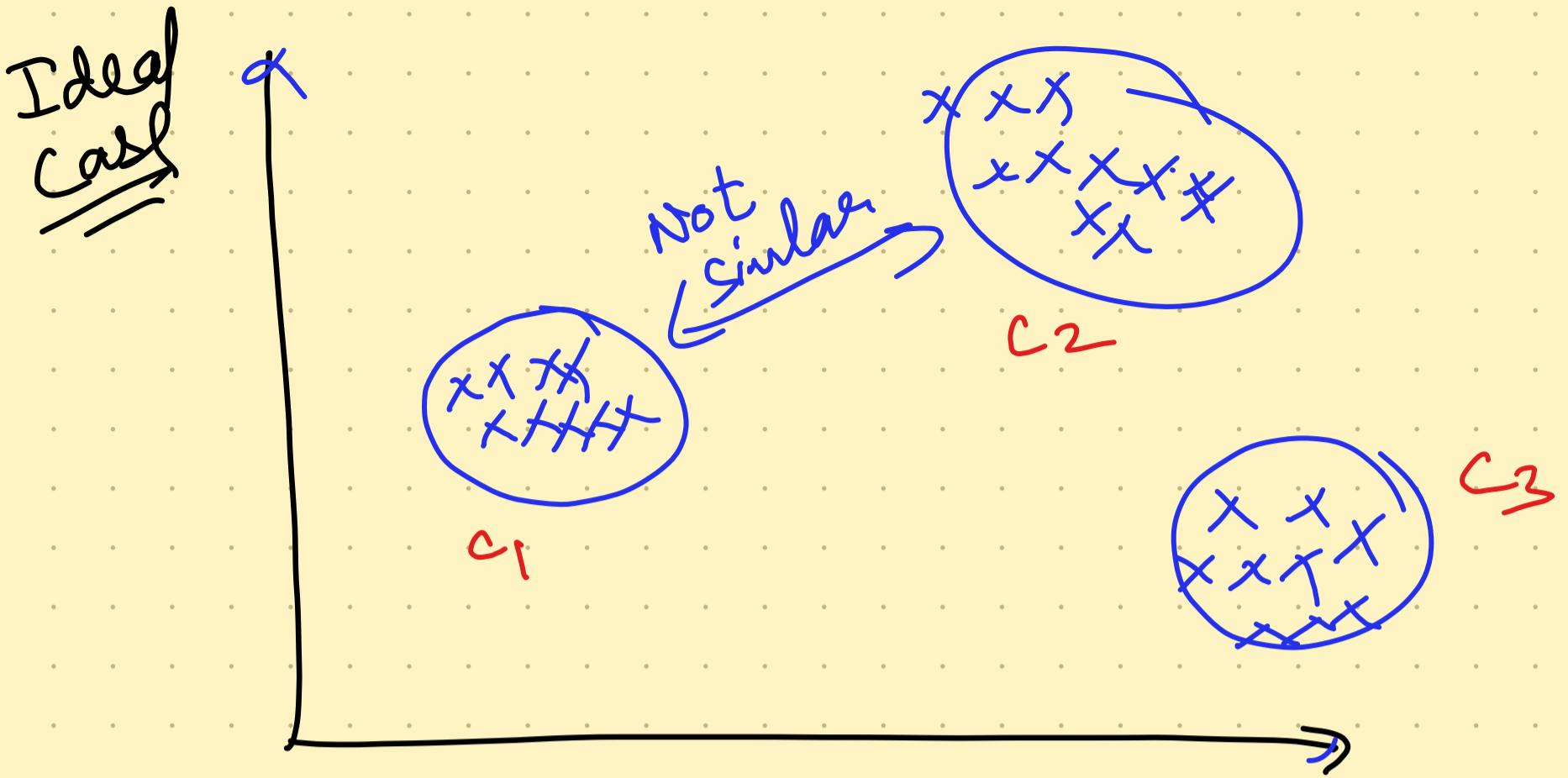




Intuition

Real world





Can we say ① Points in a cluster are close to each other

② Points in different cluster are far from each other

wicket(s)	Catch(es)
3	5
2	6
3	6



Plotted
on
paper

Q: what is the defin" of similarity

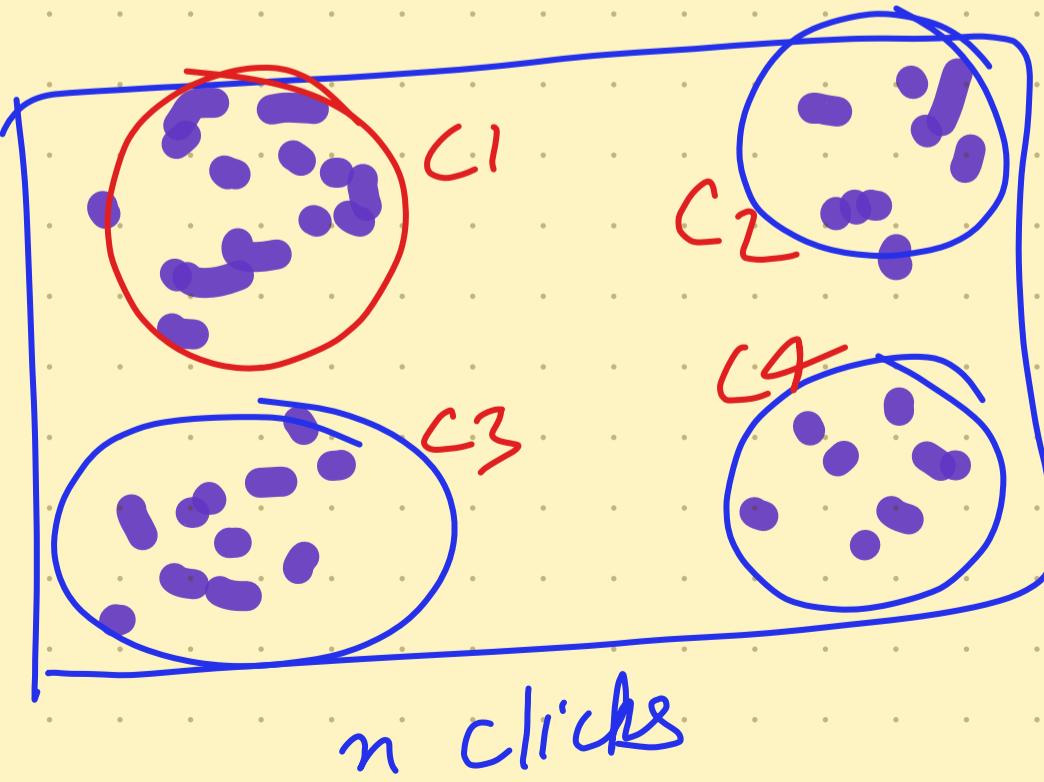
→ close pts are similar

(closeness should make Business sense as well)

ID	n_clicks	n_visits	amount_spent	amount_discount	days_since_registration
0	1476	130	65	213.905831	31.600751
1	1535	543	46	639.223004	5.689175
2	1807	520	102	1157.402763	844.321606
3	1727	702	83	1195.903634	850.041757
4	1324	221	84	180.754616	64.283300

customers spending equal amt
are similar

Amount



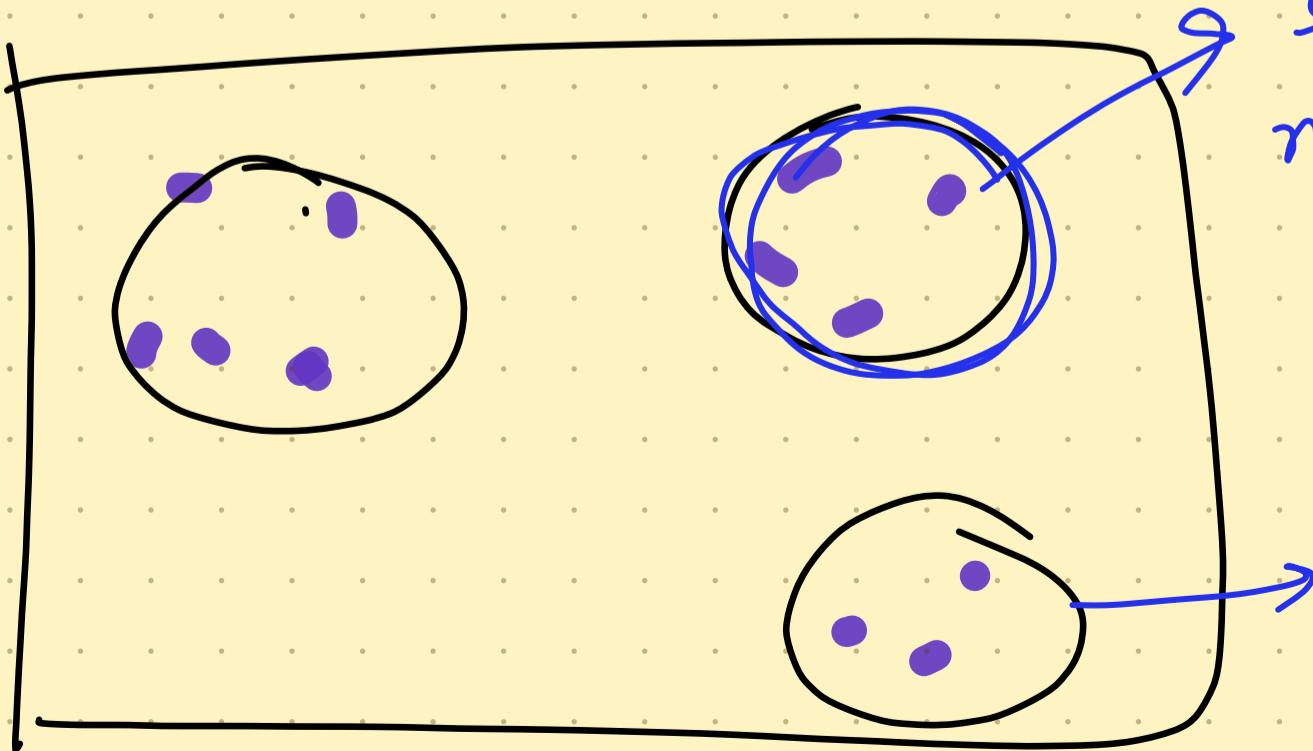
(Business decision based on the cluster/group)

n -clicks	Amount	Good/Bad

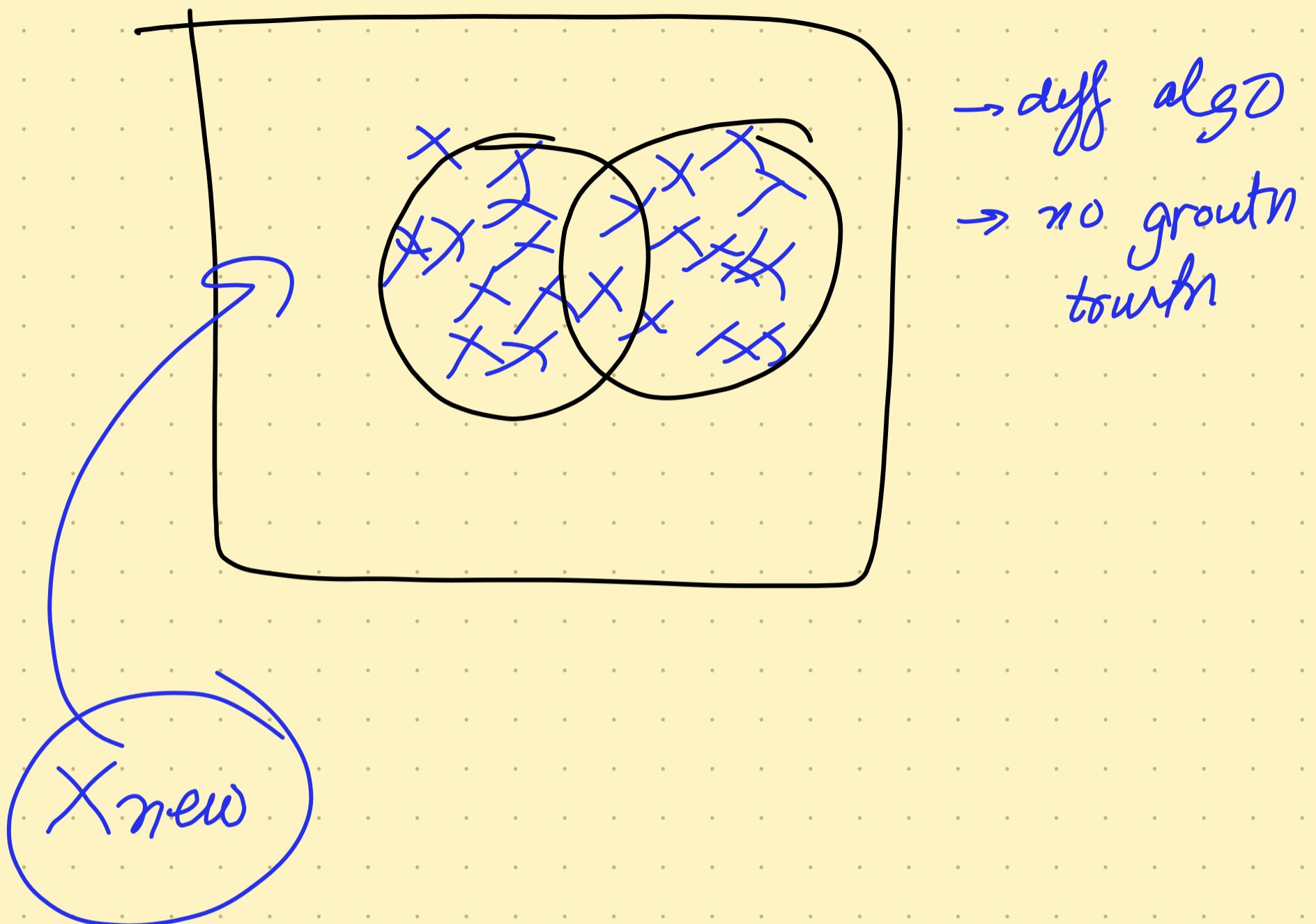
How to ensure clustering is good or Bad?

→ Ans: No one knows its good / Bad
 (we don't have Y label)

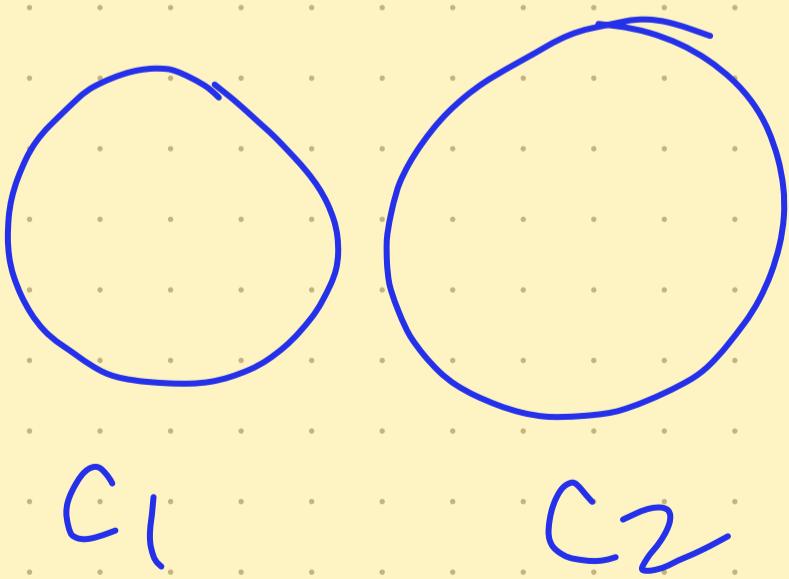
→ Business Sense

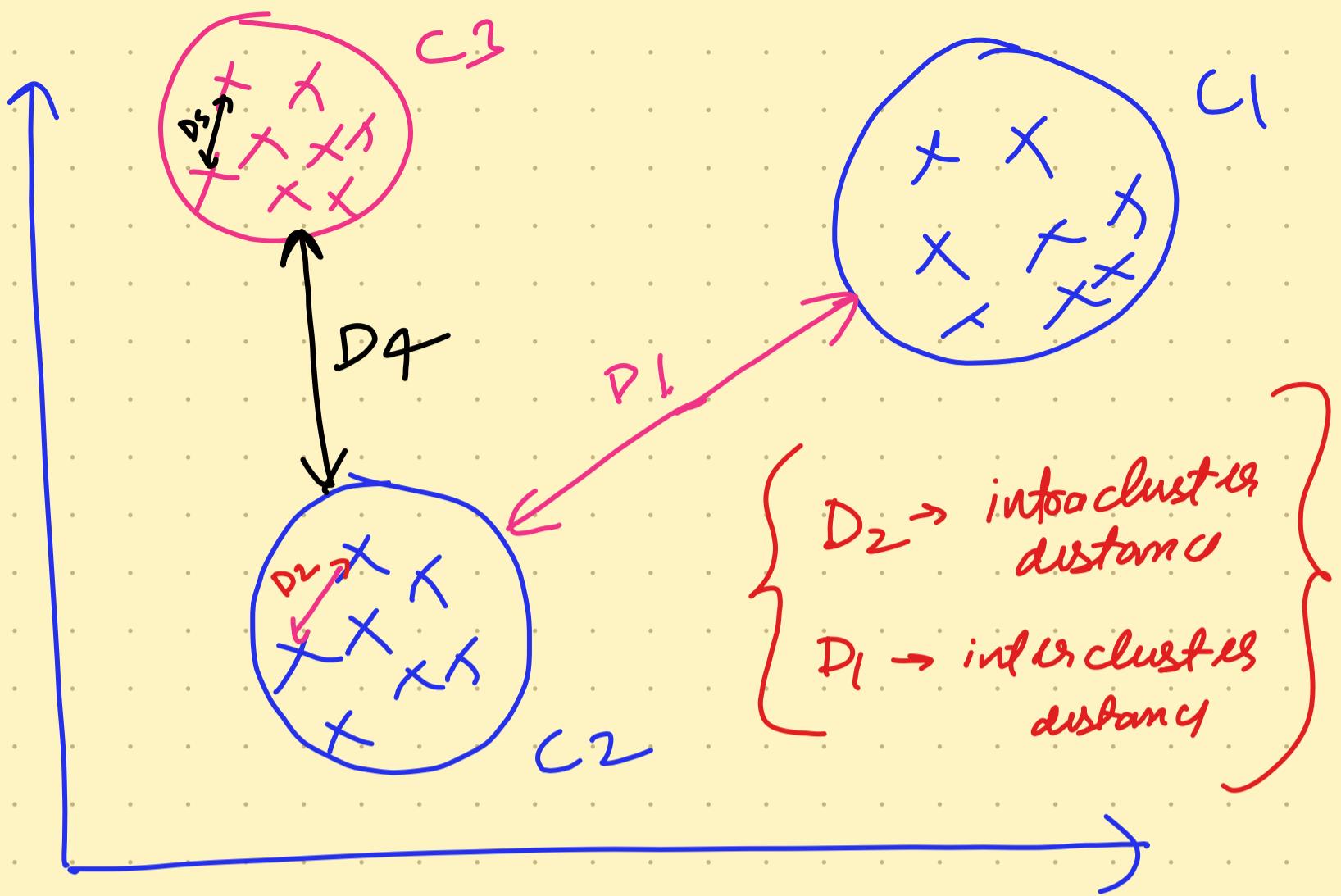


see train
 mean
 purchase
 value
 (dat analysis)



Terms and Glossary



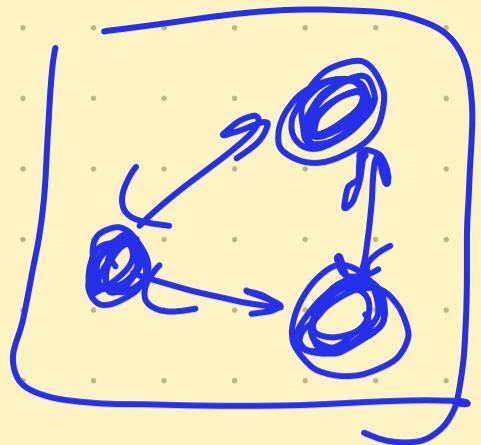


Inter College → B/w diff colleges

Intra College → within the same college

Distance

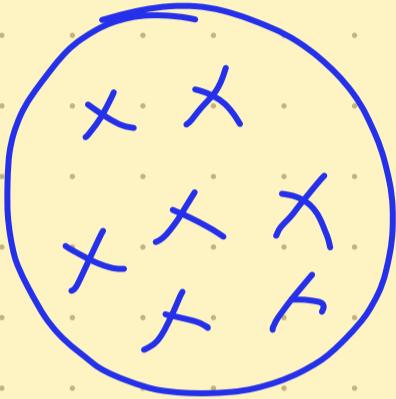
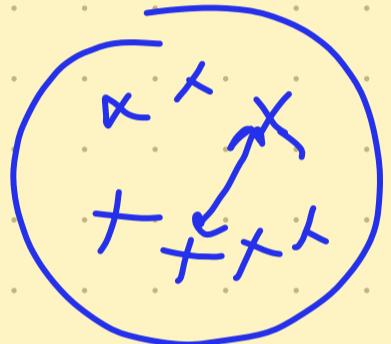
(maxim) ← Intra Cluster



Intra clusters

(minim) ←

z



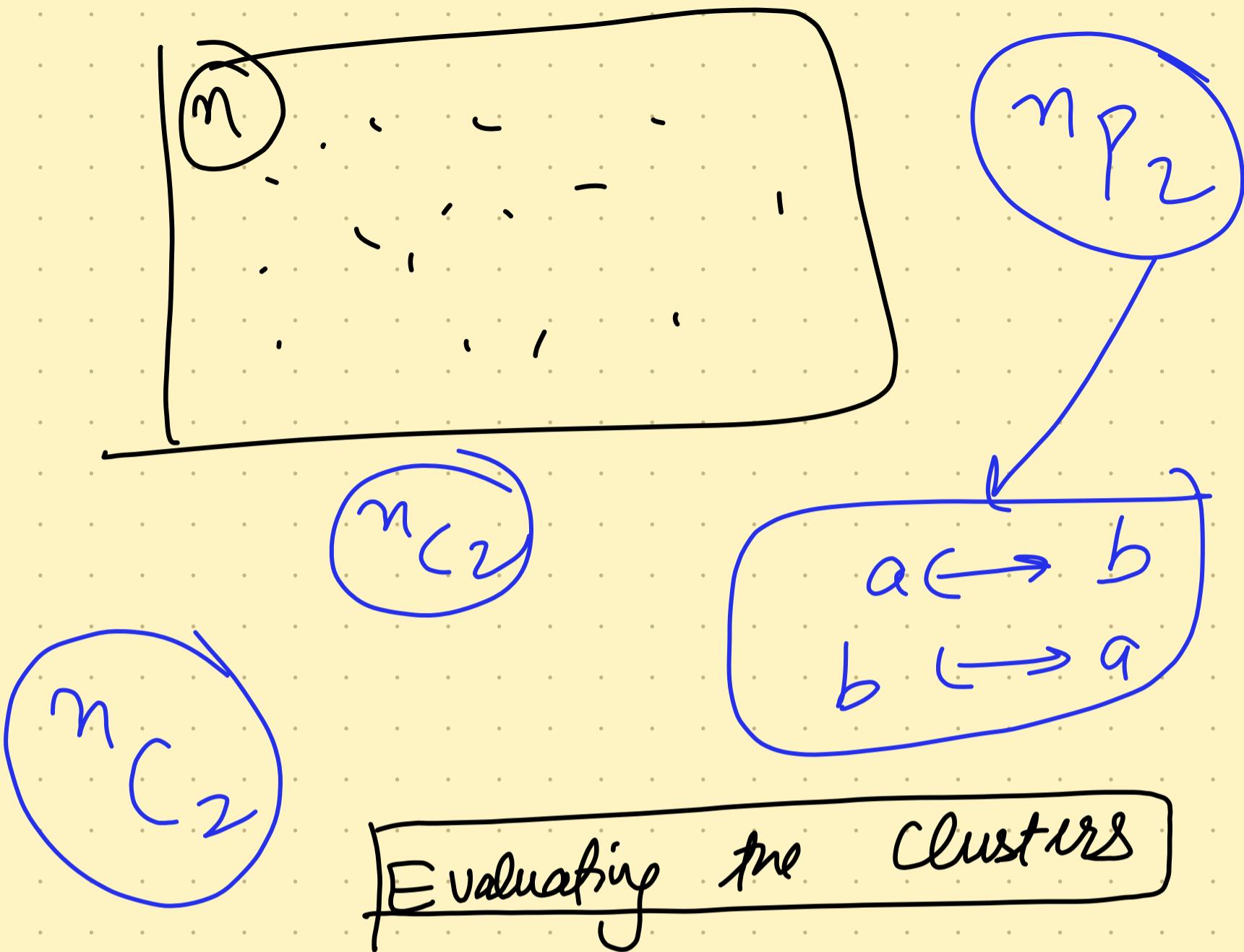
Practical

Distance Metric

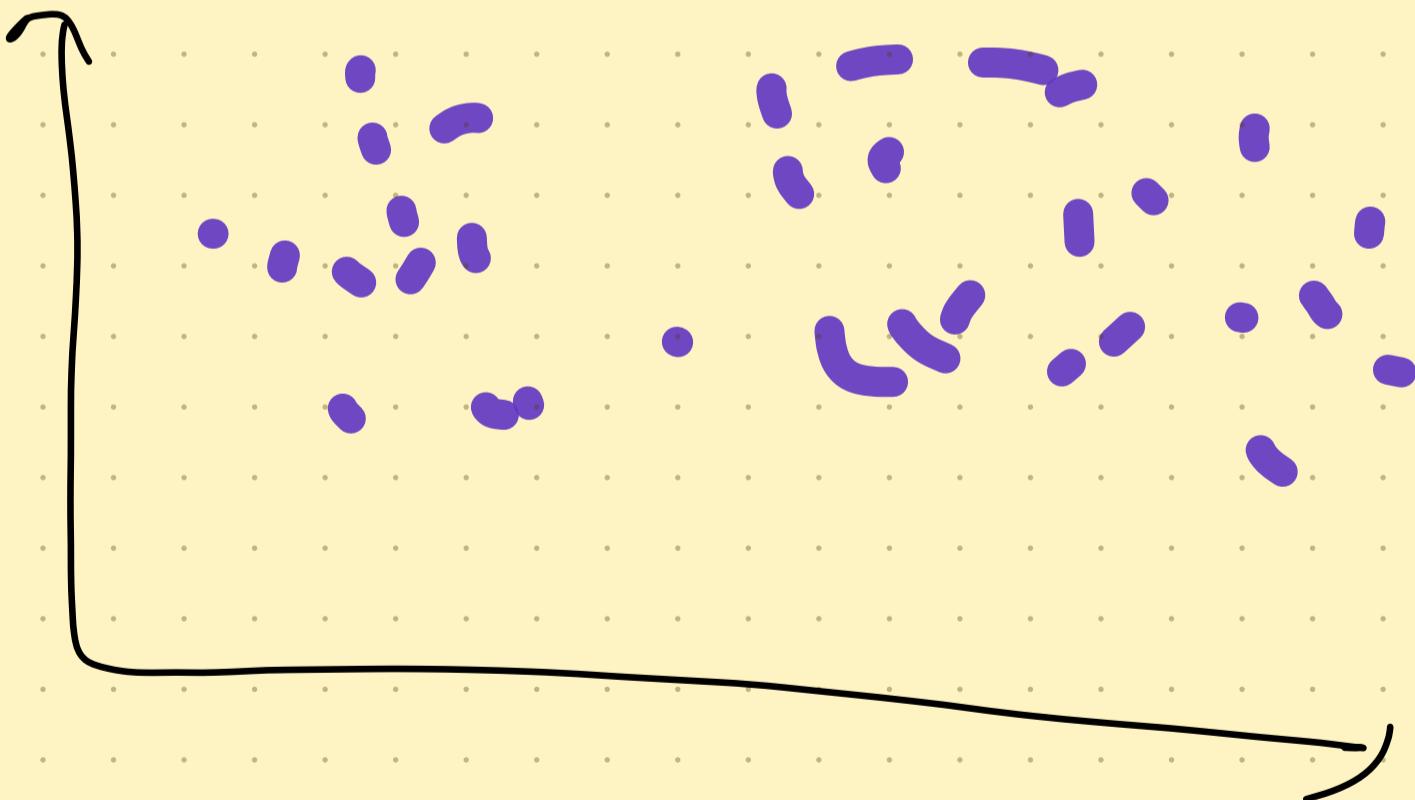
- Euclidean \Rightarrow (low dim.)
- Manhattan \Rightarrow low to med dimensions
- Cosine \Rightarrow high dimensions

Euclidean \Rightarrow

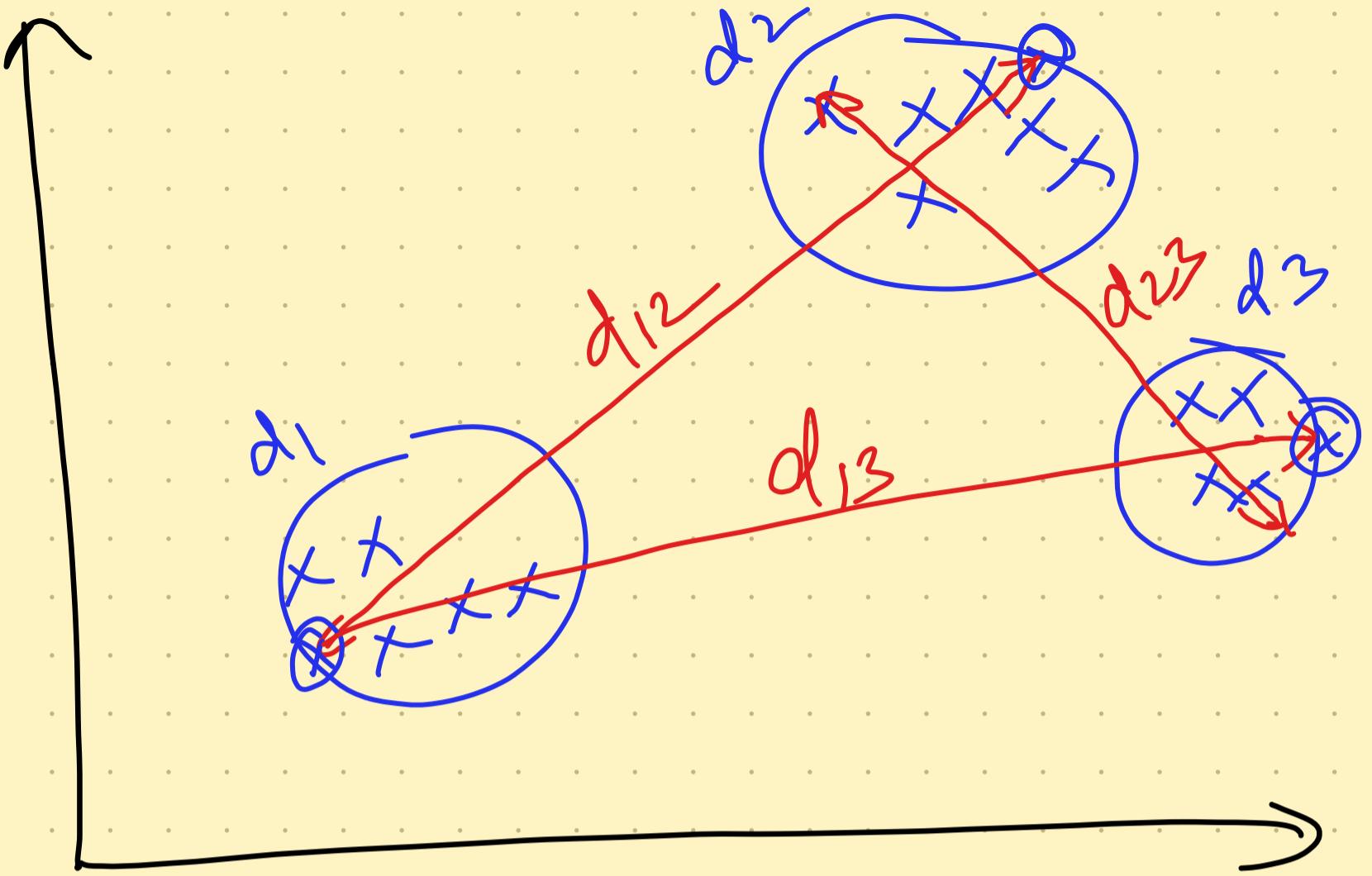
$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$



Evaluating the clusters

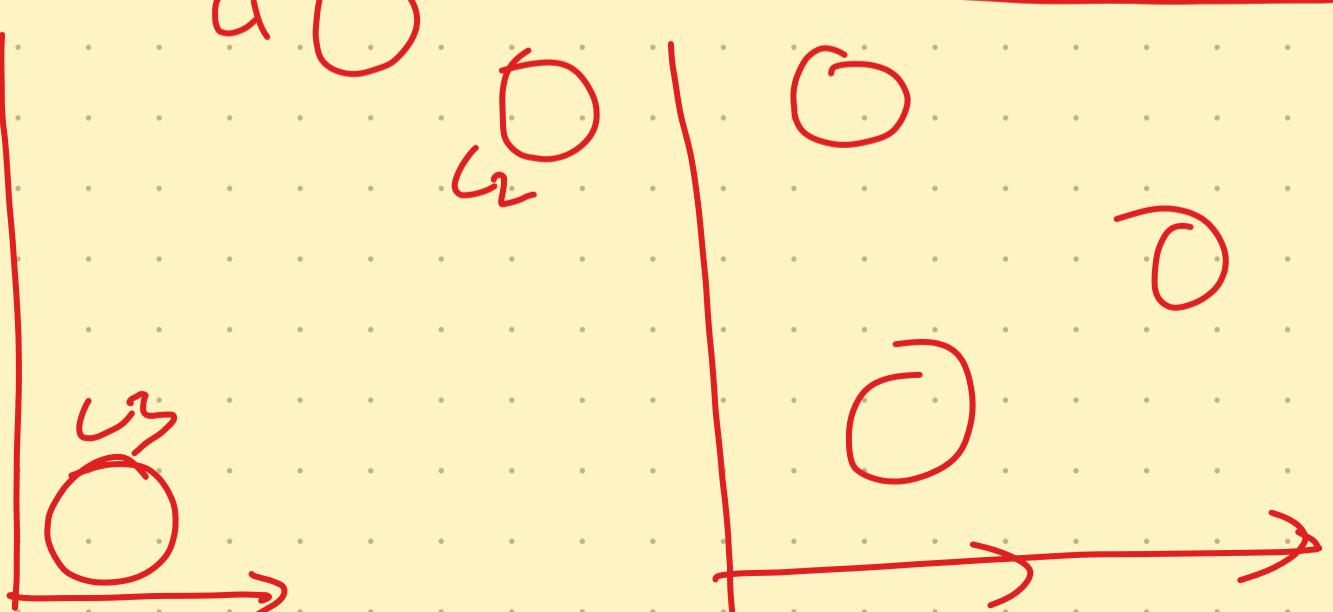


(#) [we still cannot evaluate how good the clusters are]

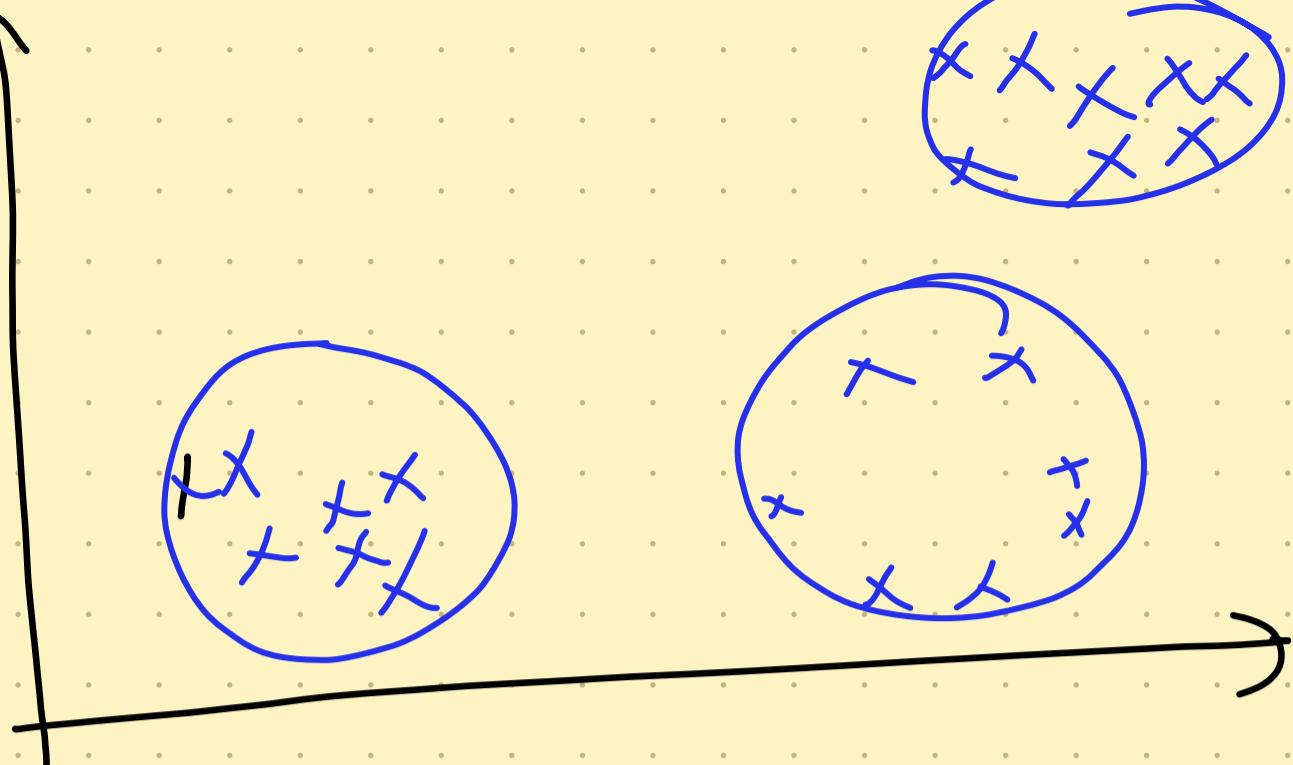


$d(i, j) = \text{distance b/w farthest points in } C_i \text{ and } C_j$ [$c_i = c_1 / c_2 / c_3$]

$$\min(d_{12}, d_{13}, d_{23}) \quad \max(d_{12}, d_{13}, d_{23})$$

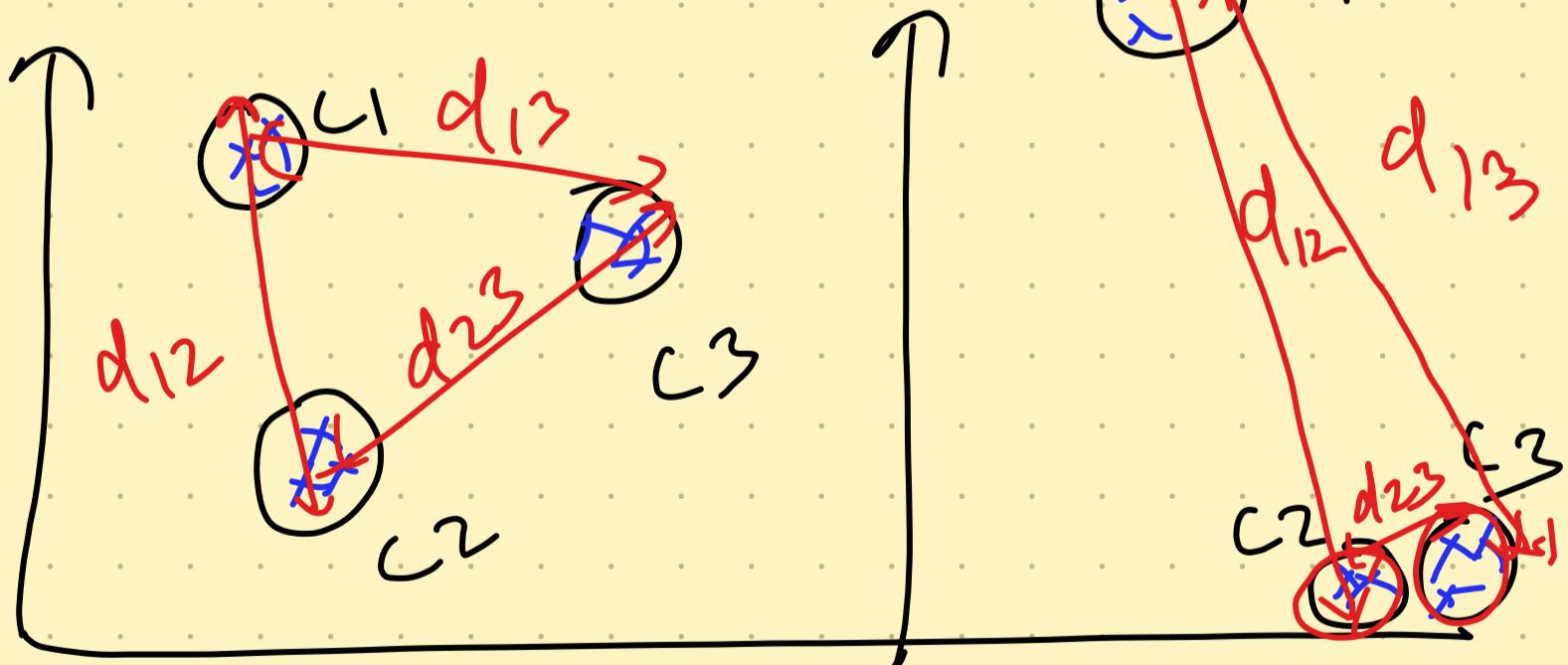


Intra cluster distance



$d'(i) = \max$ distance b/w 2 points in C_i

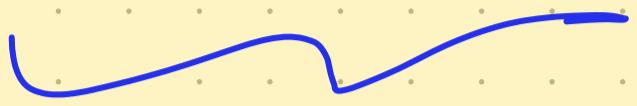
Q which is good cluster



$\max(d_{12}, d_{13}, d_{23}) \Rightarrow h$

$$\min(d_{12}, d_{13}, d_{23}) = \gamma$$

$$\gamma_{\text{(Case I)}} > \gamma_{\text{(Case II)}}$$



Case I is

a better cluster arrangement

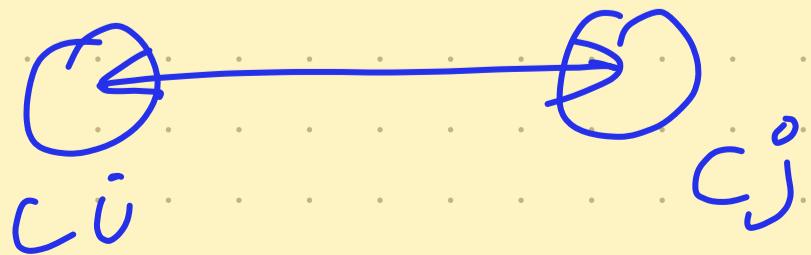
Interc Cluster Distanc

$$\text{Dunn Index} = \frac{\min(\text{Inter clusters dist})}{(\max)(\text{Intra cluster Dist})}$$

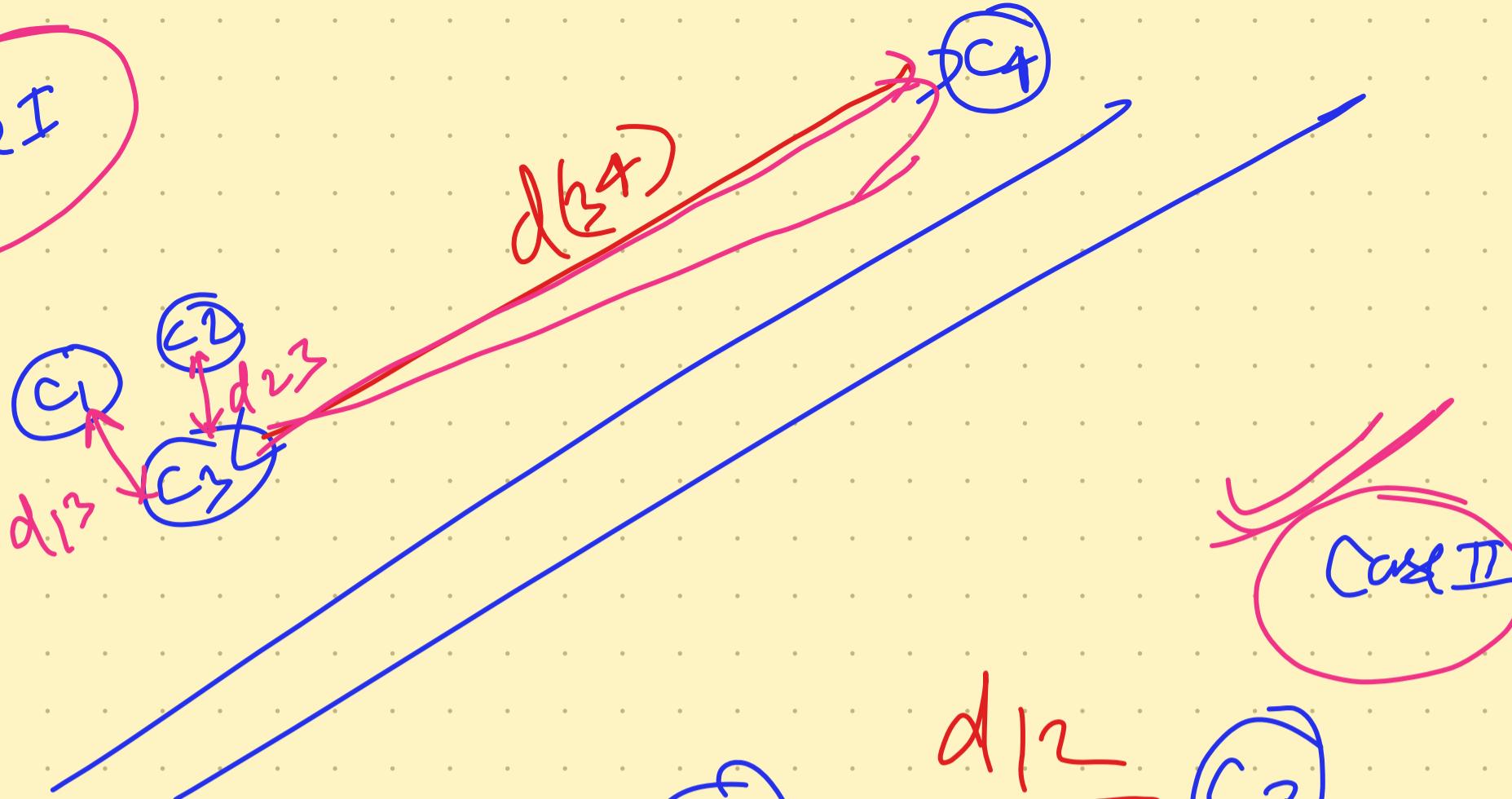
$$D \cdot I$$

$$= \frac{d(i, j)}{d'(i)}$$

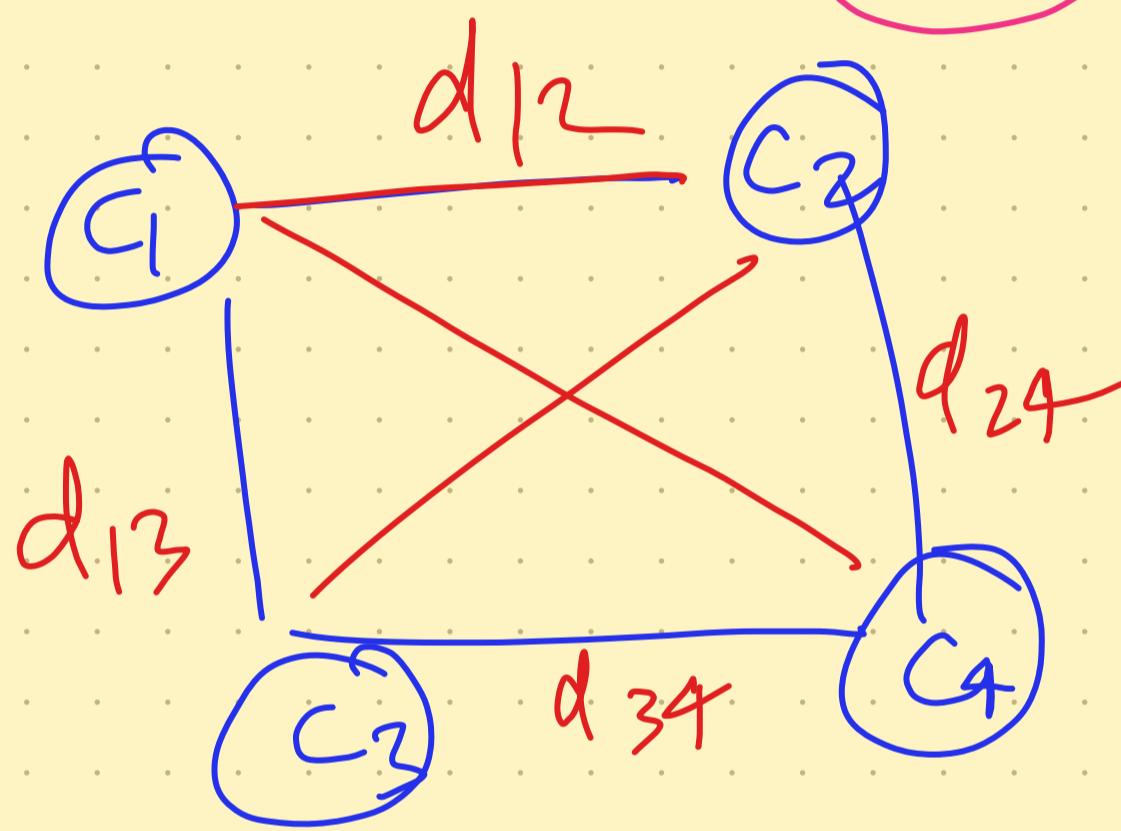
i, j are from diff clusters



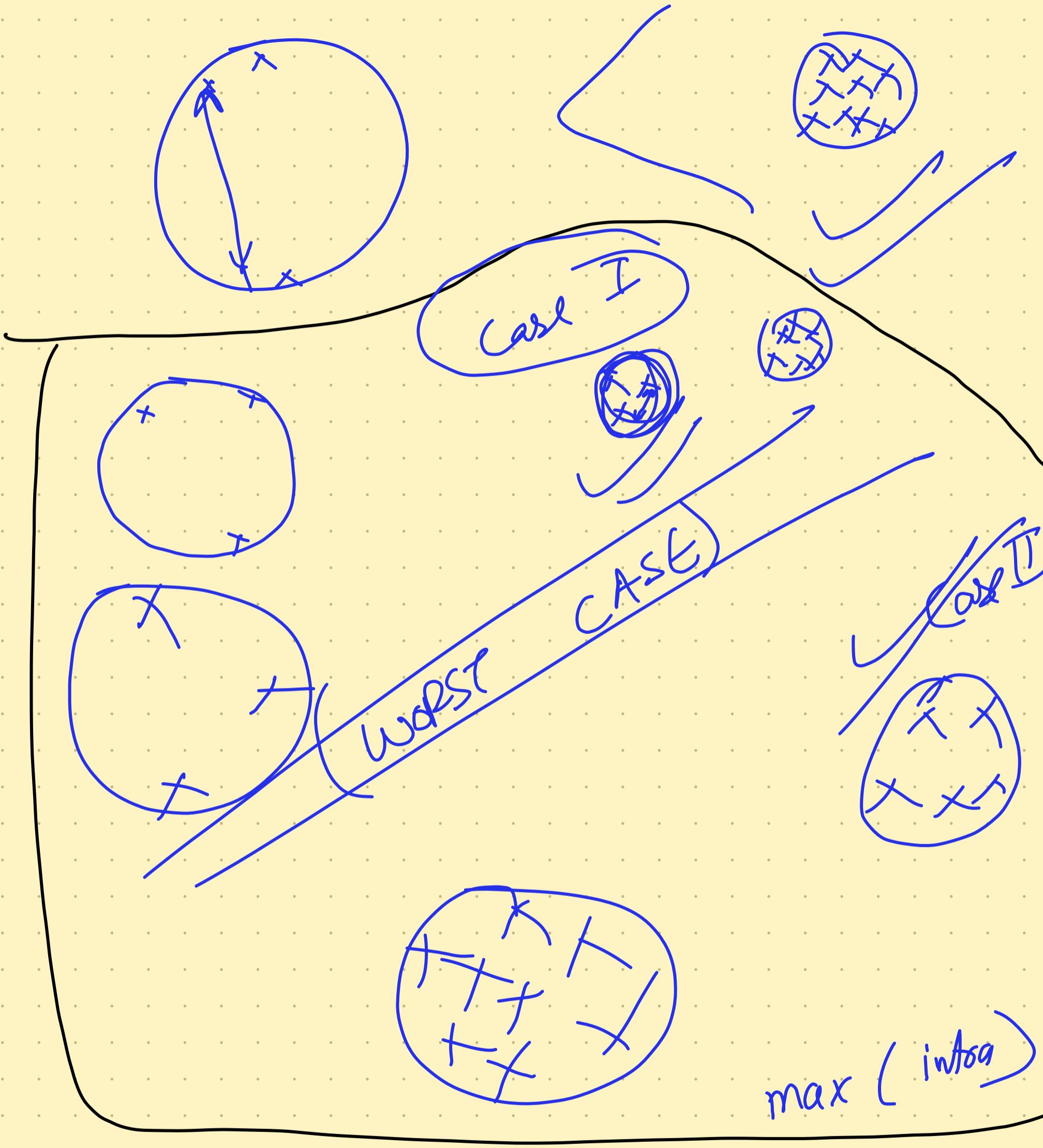
Case I



Case II



$\min \frac{\text{inter cluster distance}}{\max \text{ intra cluster}}$

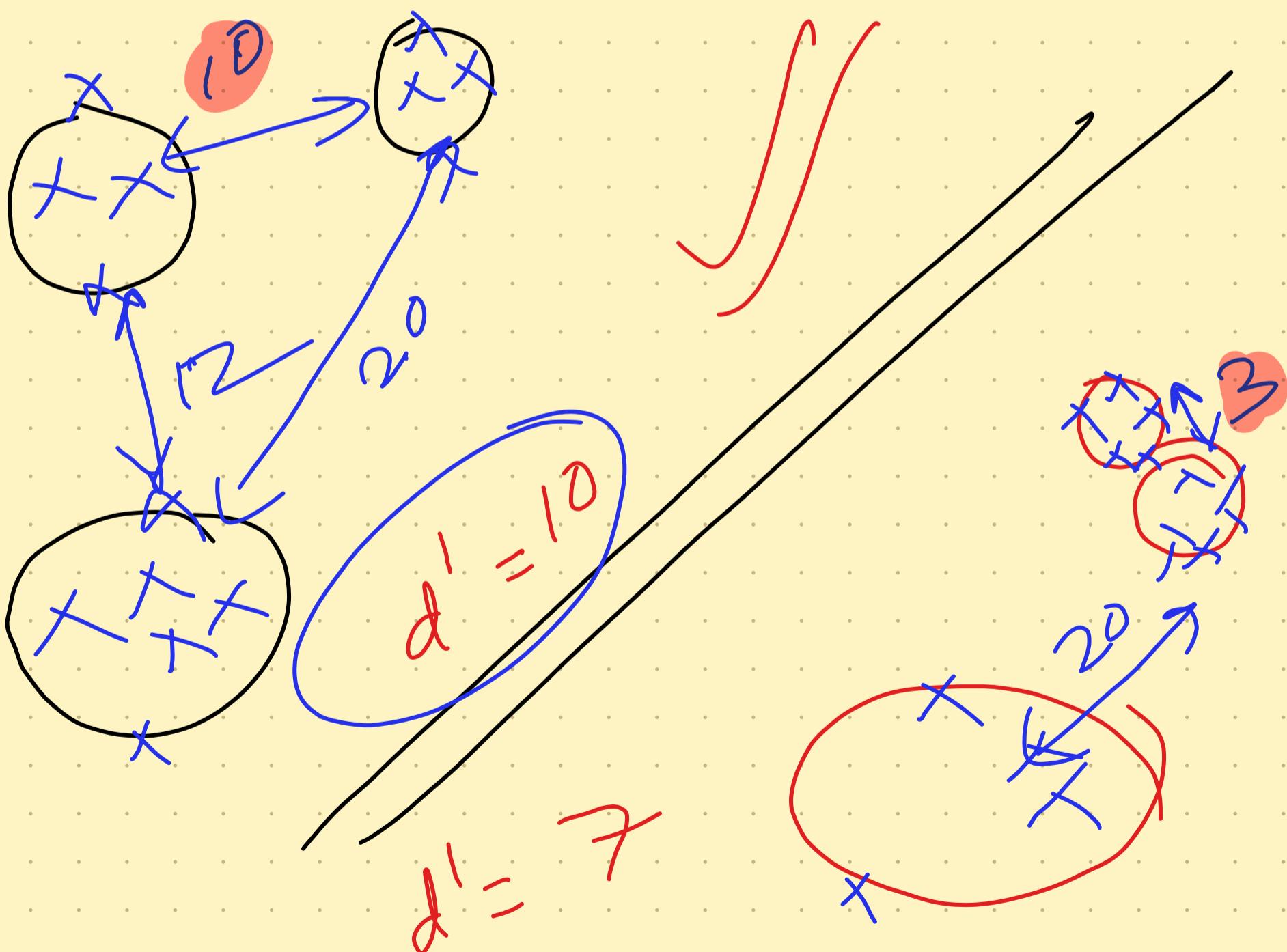


Drun Indx

$$\text{Dunn Index} = \frac{\min(\text{Inter})}{d' = \max(\text{Intra})}$$

(10)

$$\min(\text{Inter}) = 10 \text{ (max Inter)}$$



$$\text{Dunn Index} (\mathcal{D}) = \left(\frac{0}{10} \right)$$

$$D\text{-}I(\mathcal{D}) = \left(\frac{3}{3} \right)$$

Higher Dunn Index $\frac{19}{3} \Rightarrow$ Better clustering Algo

$$O(n \log n)$$

Applications

Jeans

10 million images

Top

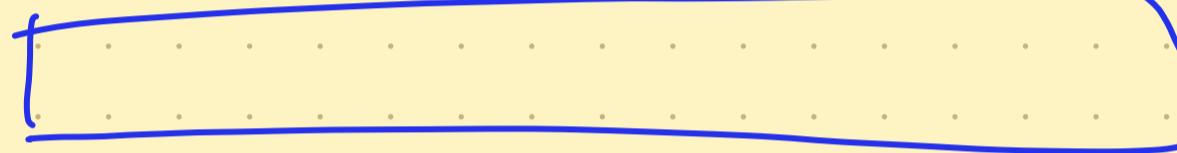
Shirt

Night
Gown

$d = 10^D$

d - dimension vector

image =

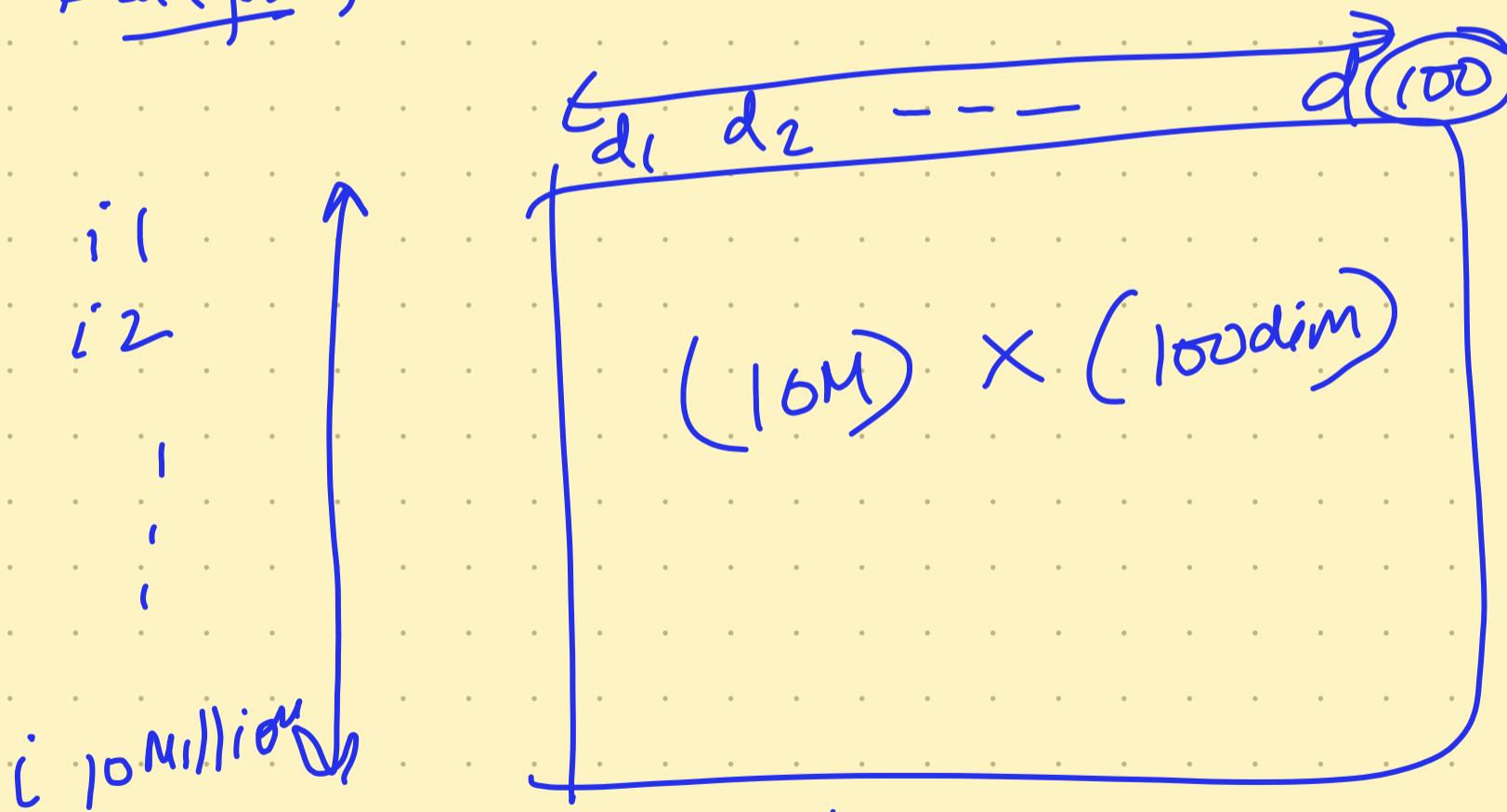


ID

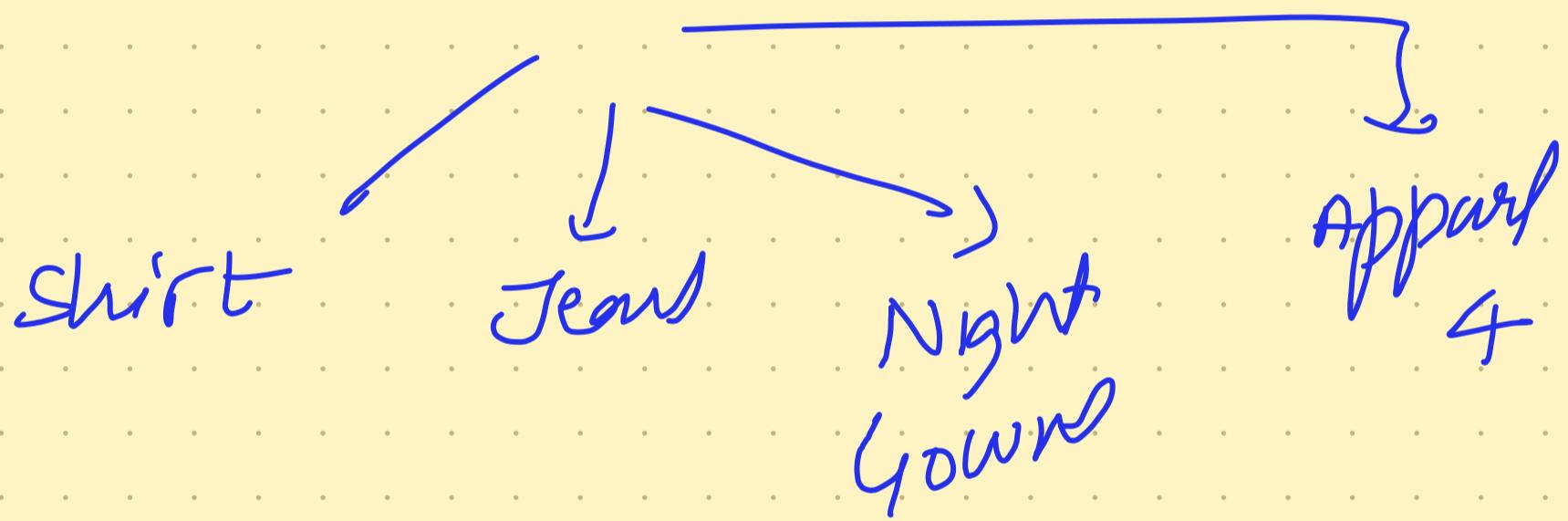
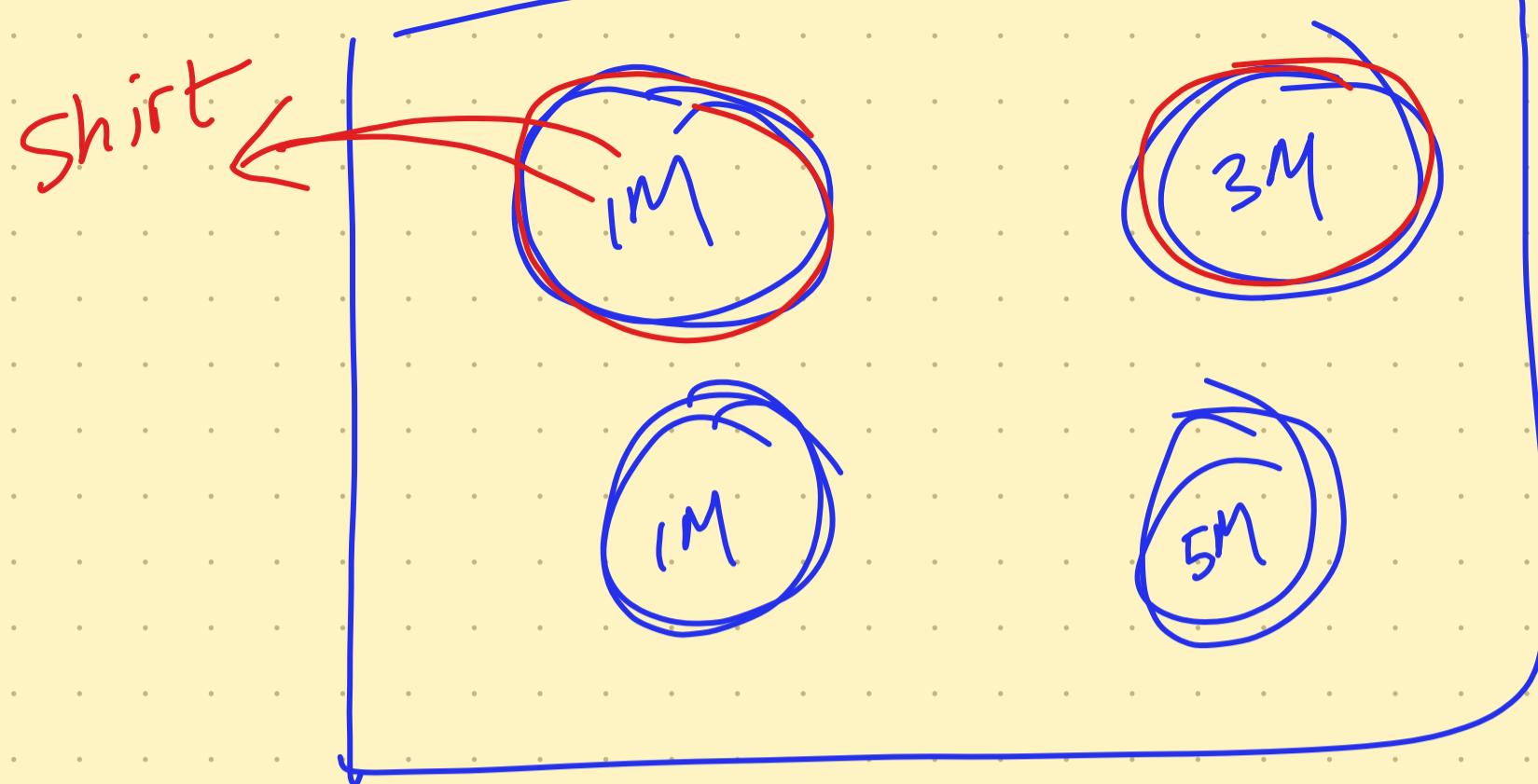
image

Data form

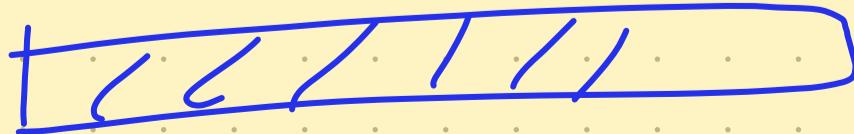
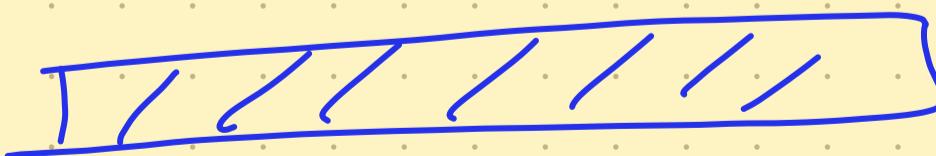
$N = 10 \text{ Million}$



clustering

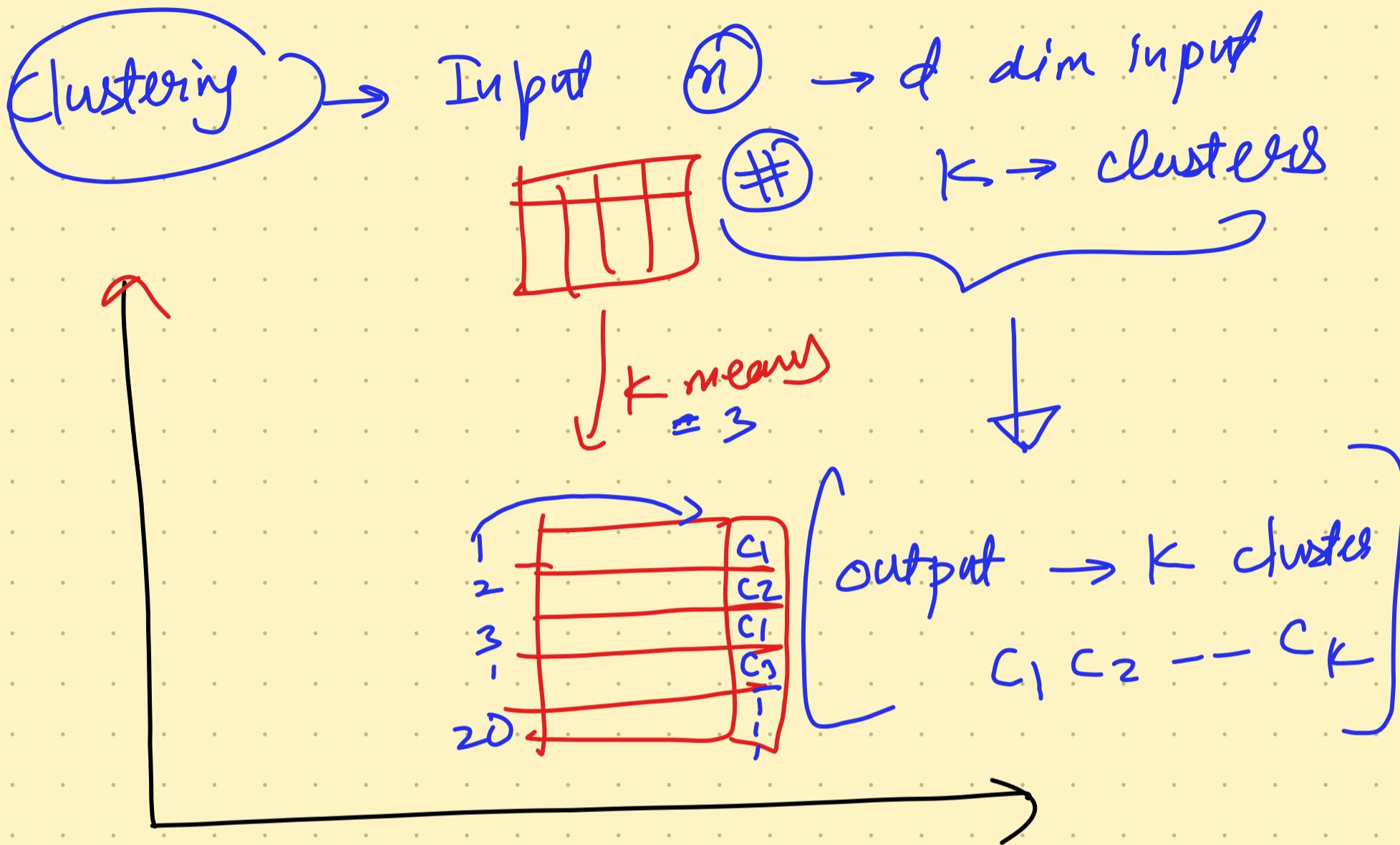


Document Labeling →



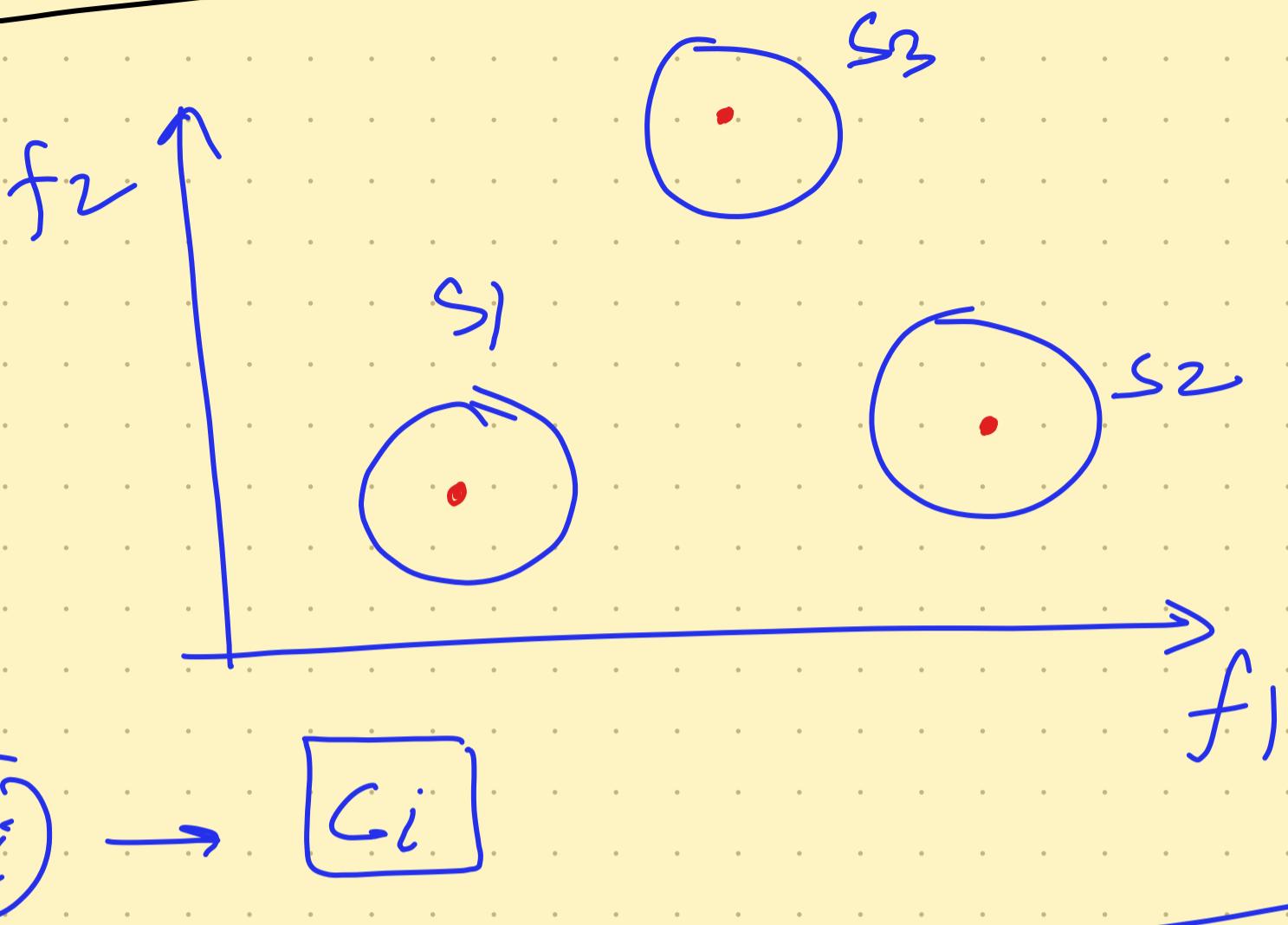
K Means Algorithm

$K \rightarrow$ no of clusters



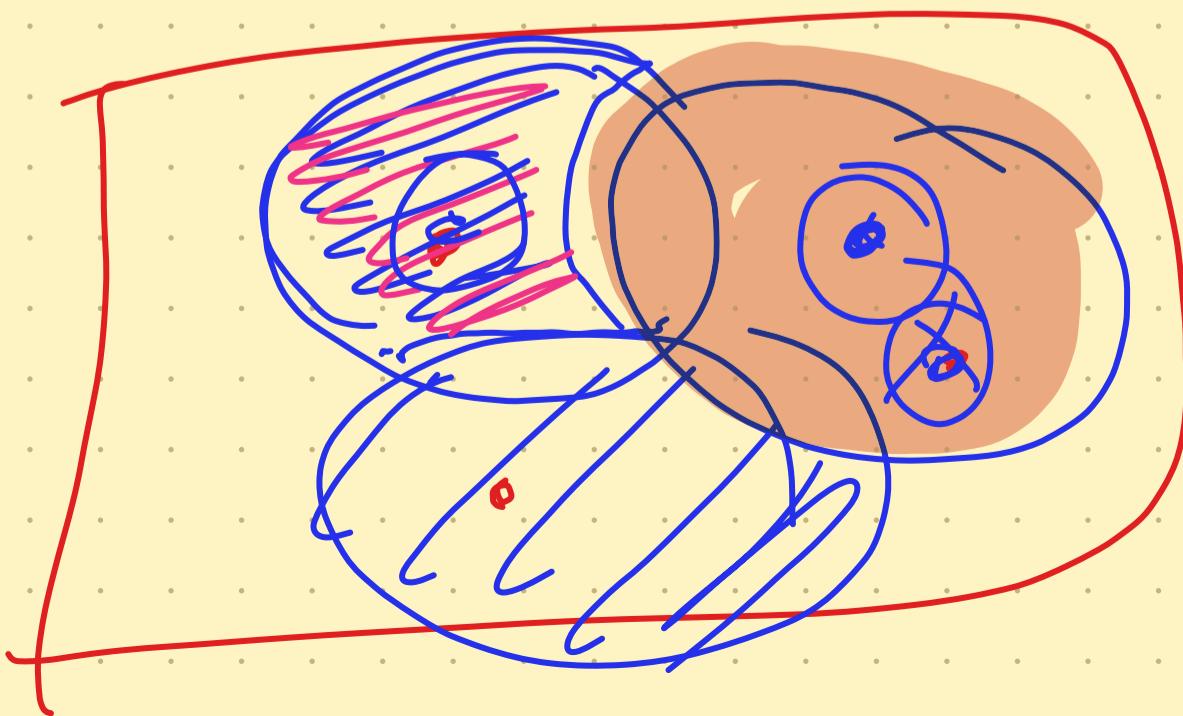
- K means
- simple
 - optimzⁿ Algo
 - Popular
 - you need to define K before hand

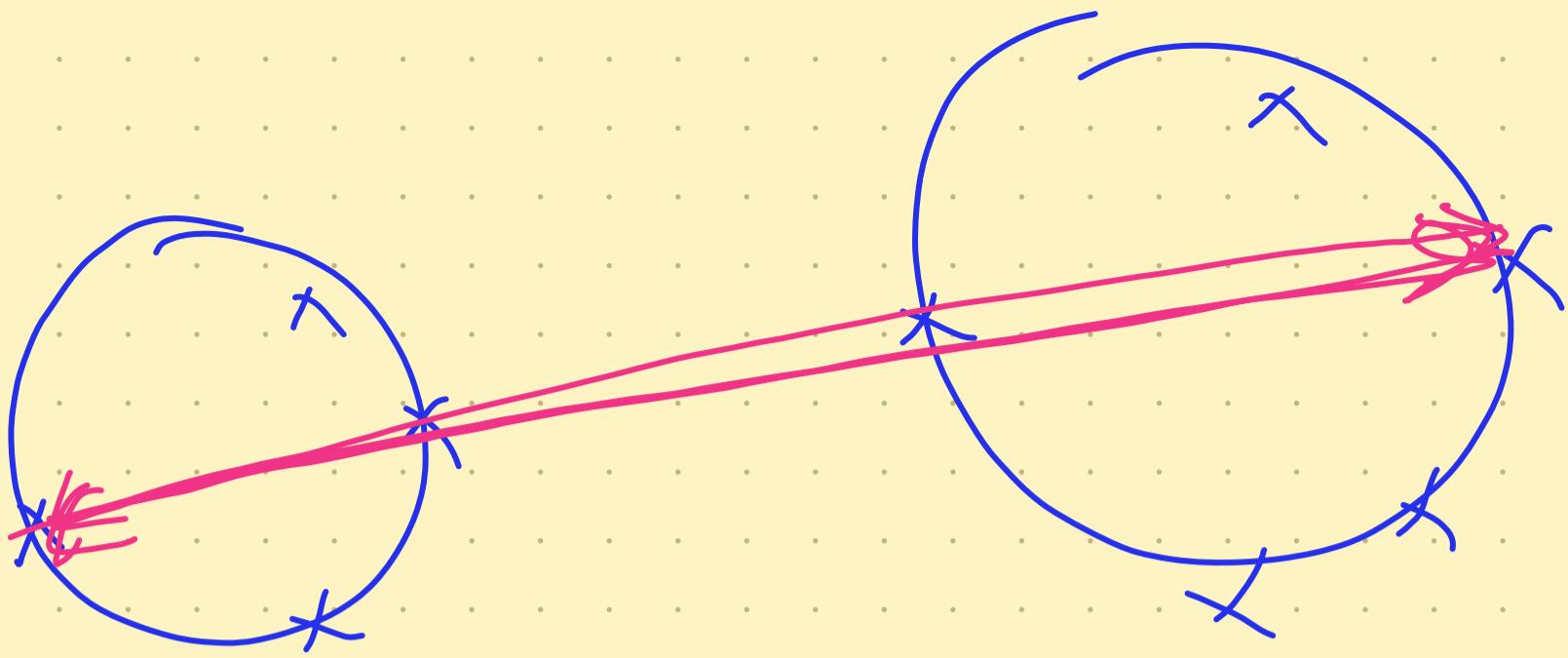
2 Dim dataset \rightarrow 3 clusters ($K=3$)



$$S_i \cap S_j = \emptyset \text{ and each } \pi_i \text{ gone } S_j$$

Mean Data $= \left(\frac{\pi_1 + \dots + \pi_m}{m} \right) \in S_i$





Set of parameters
(txt file)

R

Model

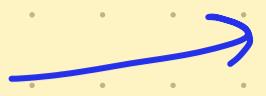
C set of parameters

Input

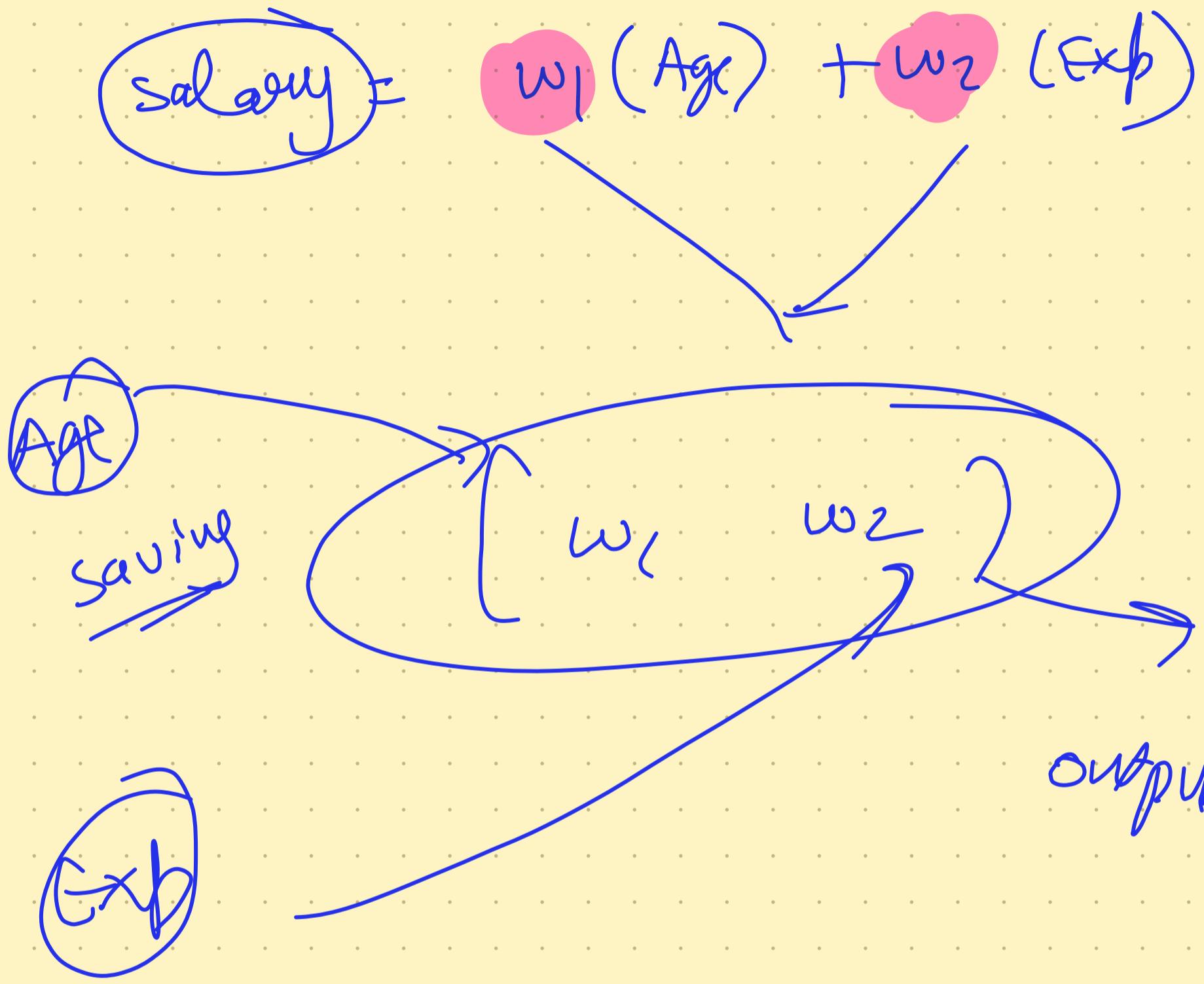
[0 0 0 0]

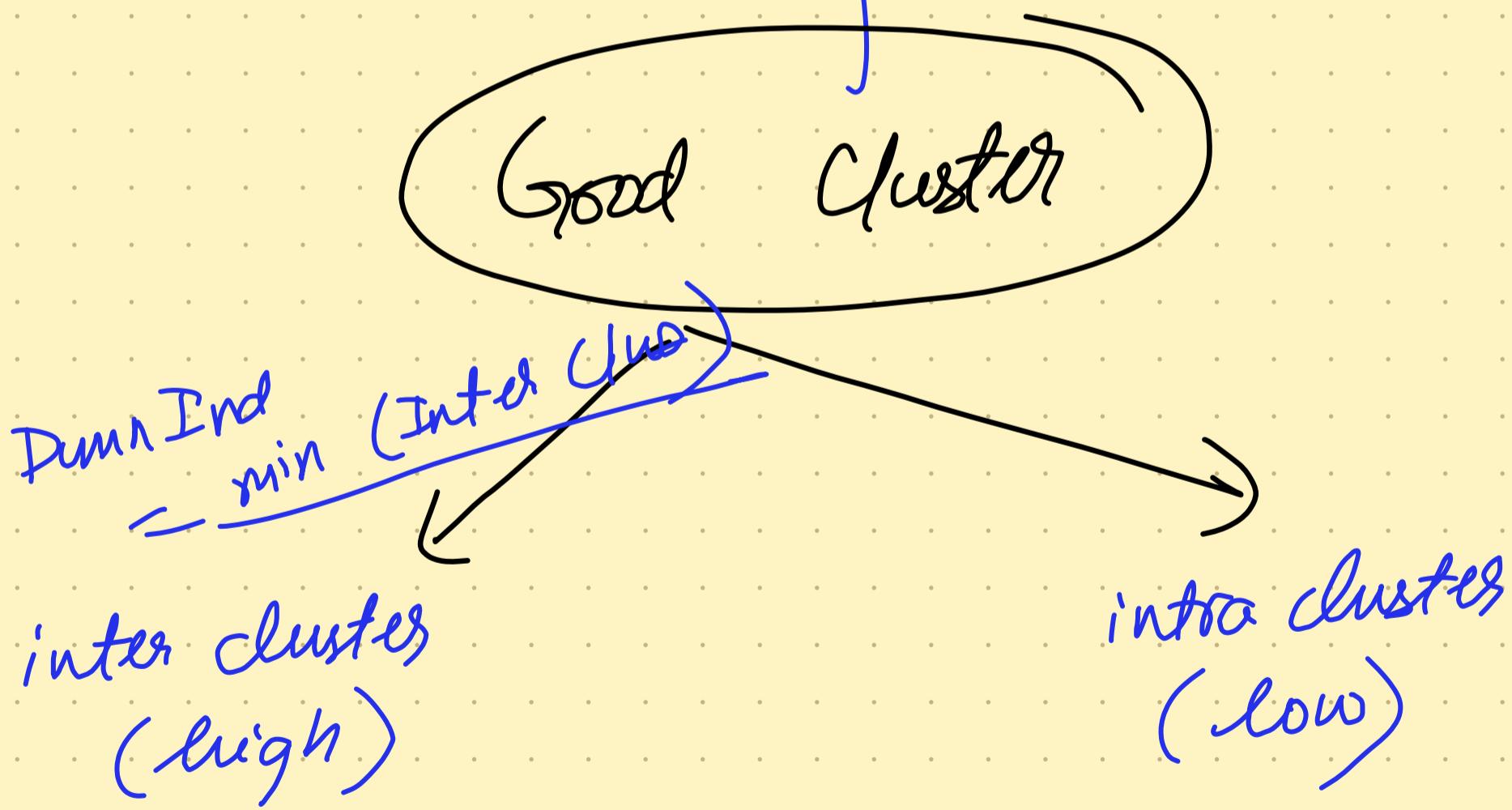
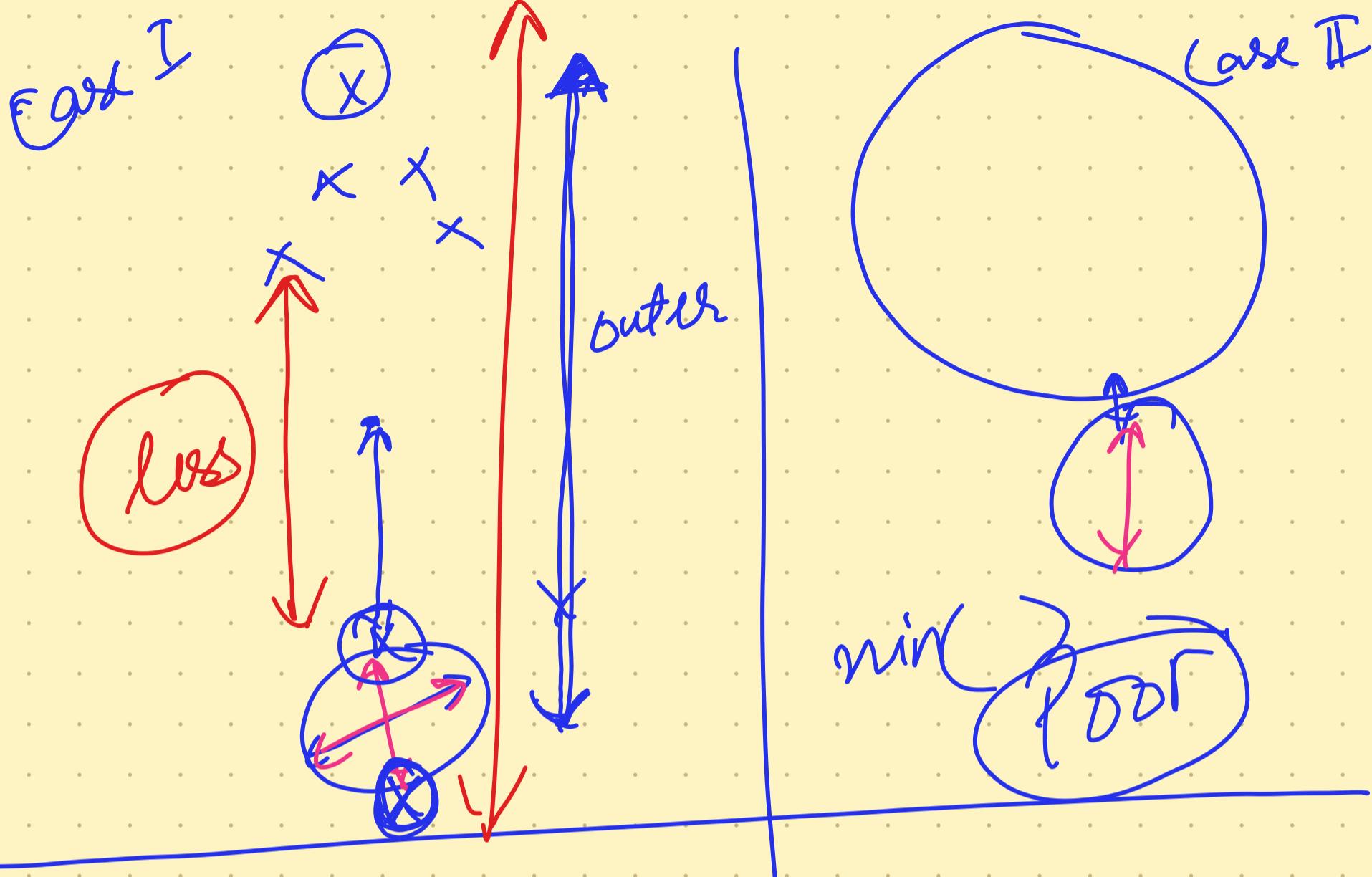
Output

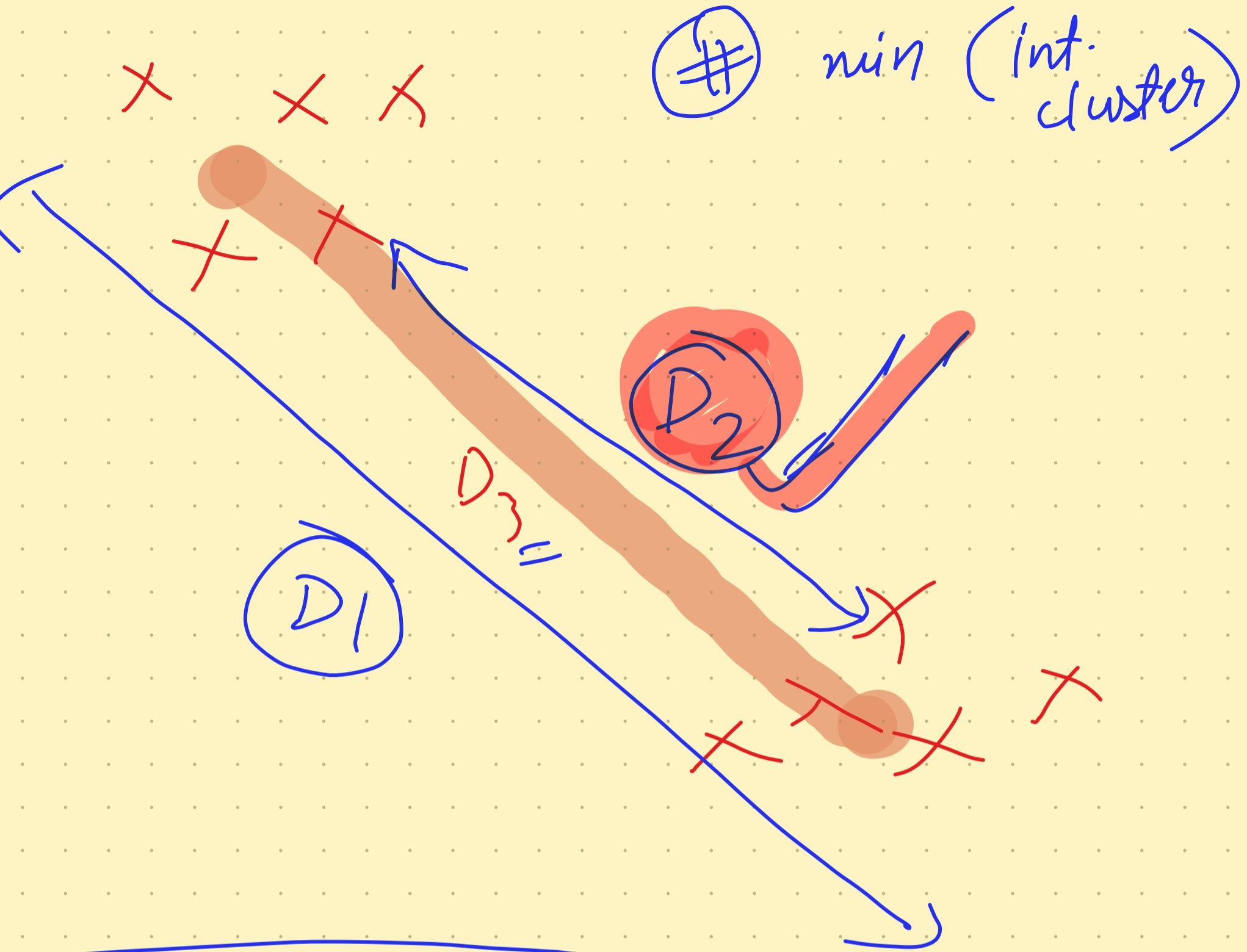
L.R



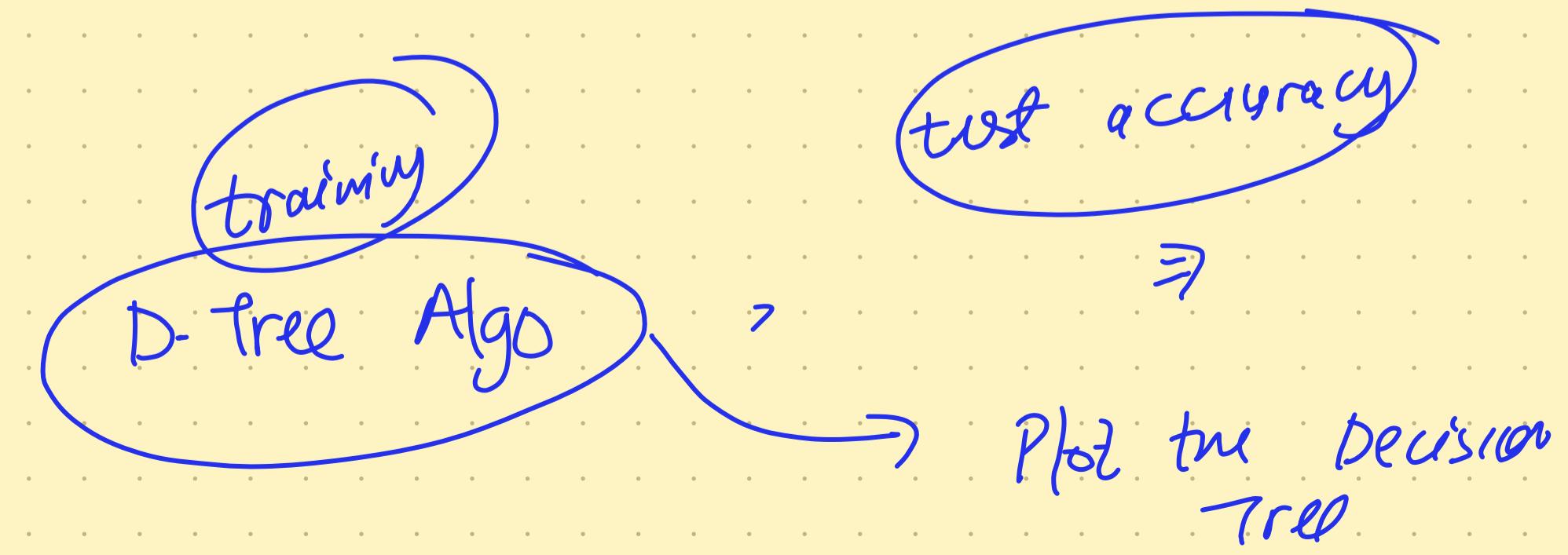
weights of model

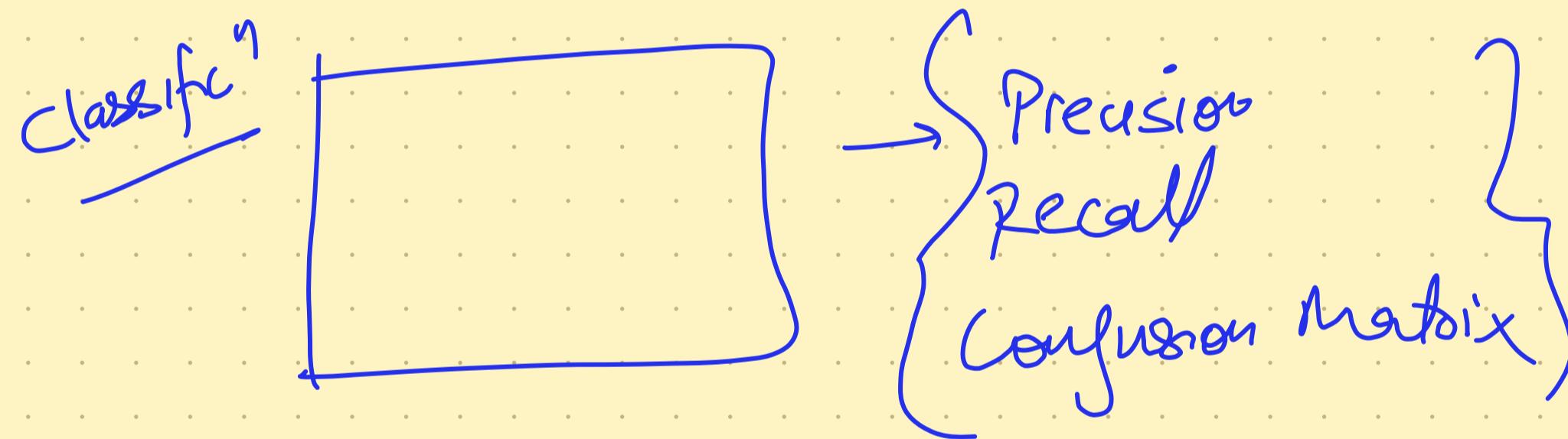
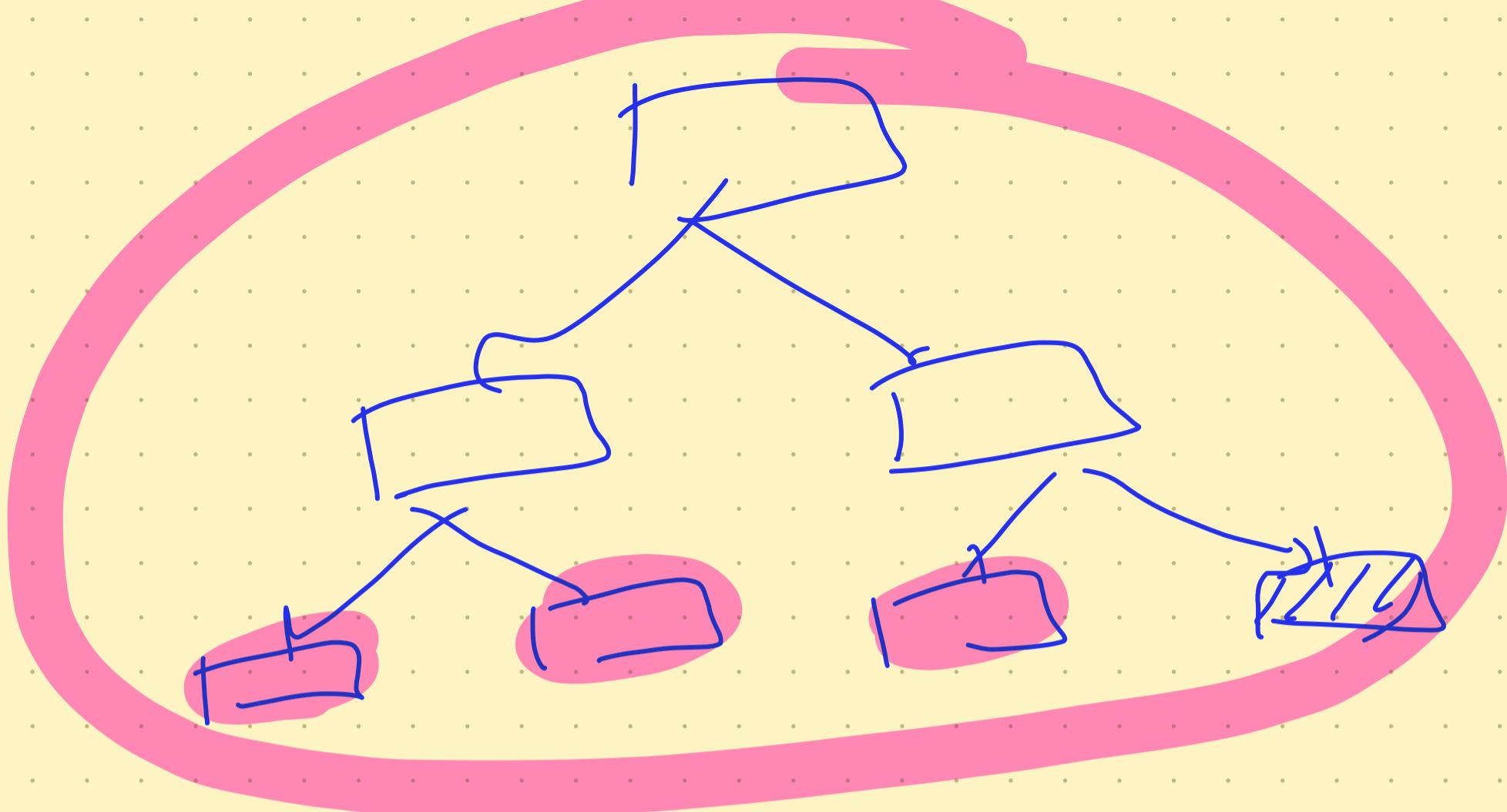






P_1 and P_2

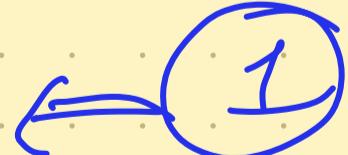


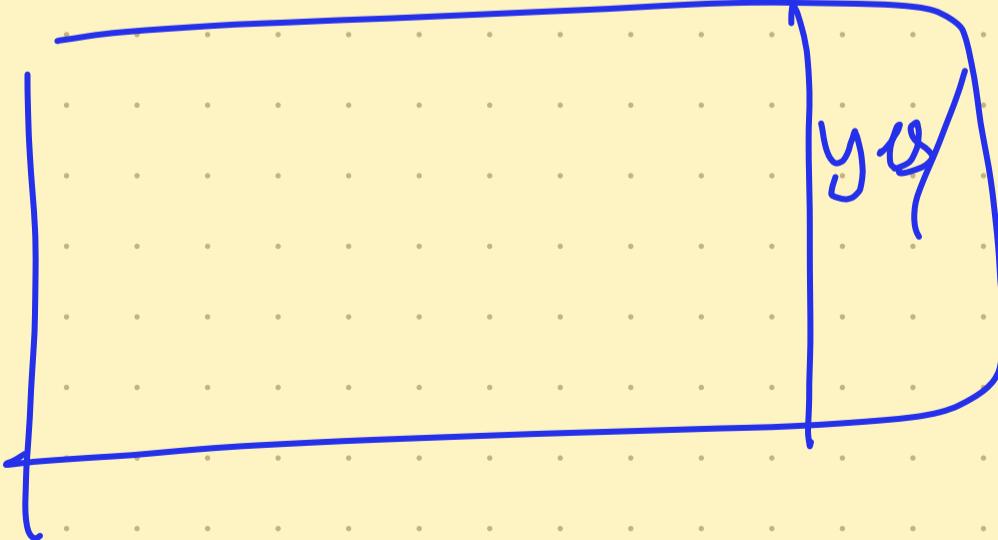
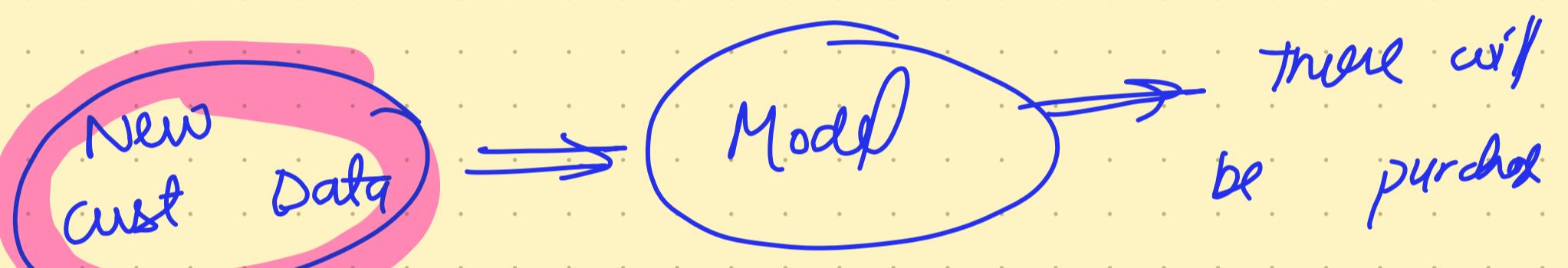
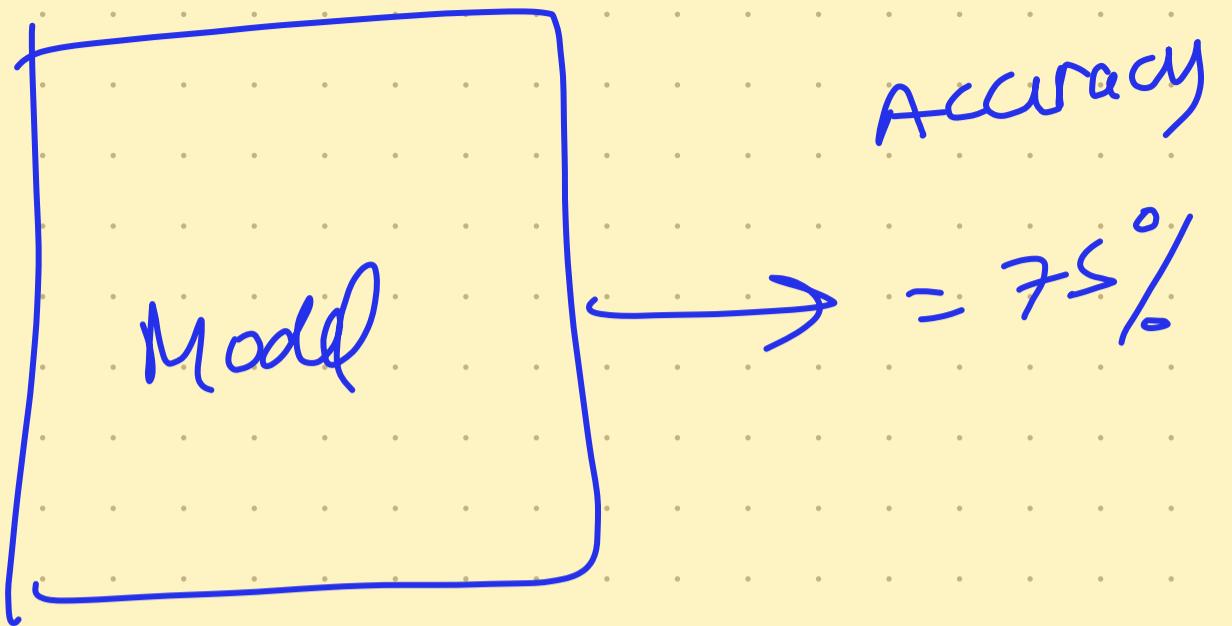


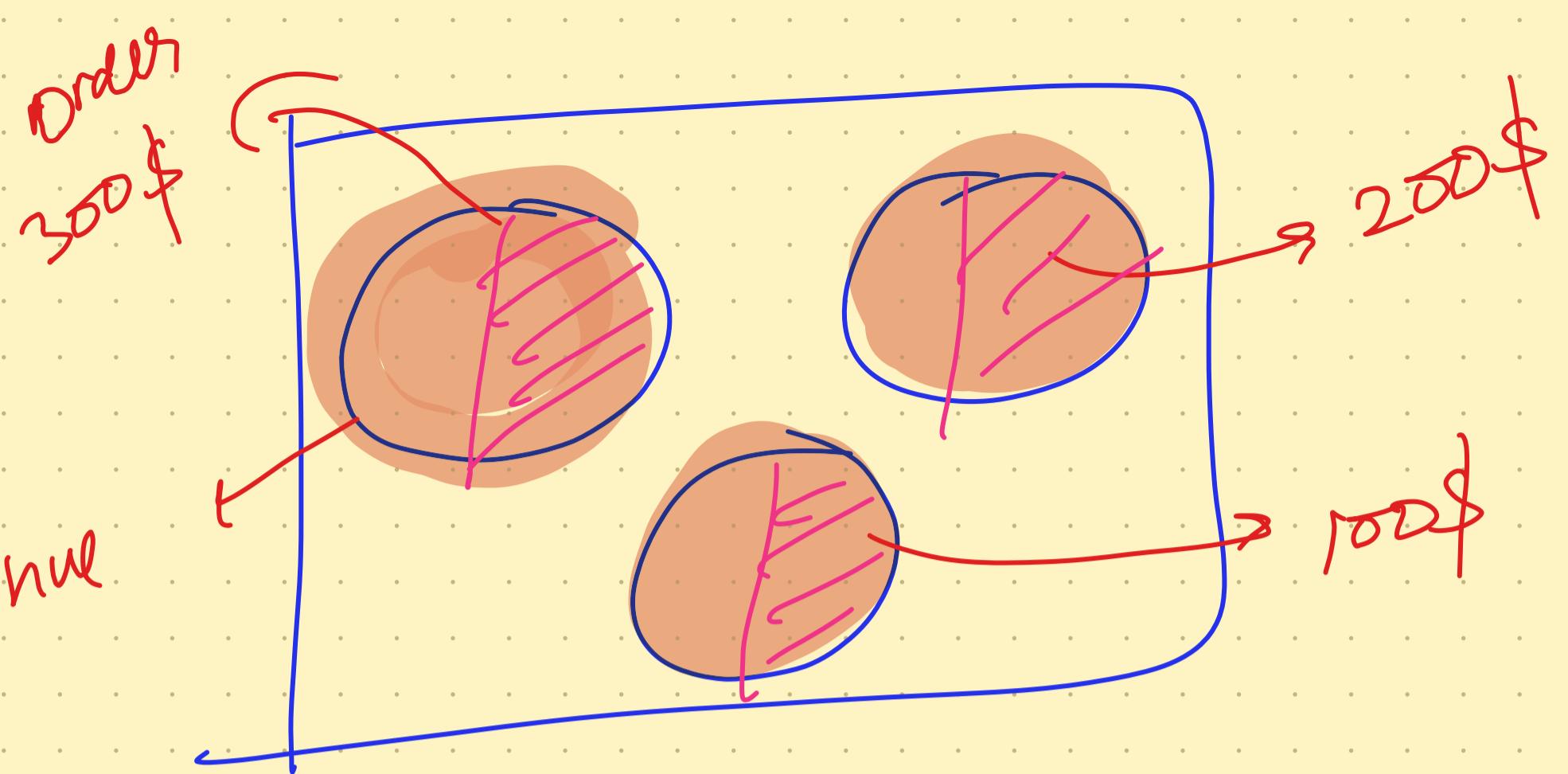
Purchas / Not → 75

Accuracy = 75

Acurr > 50%







Yes / No